

# Cell type prediction

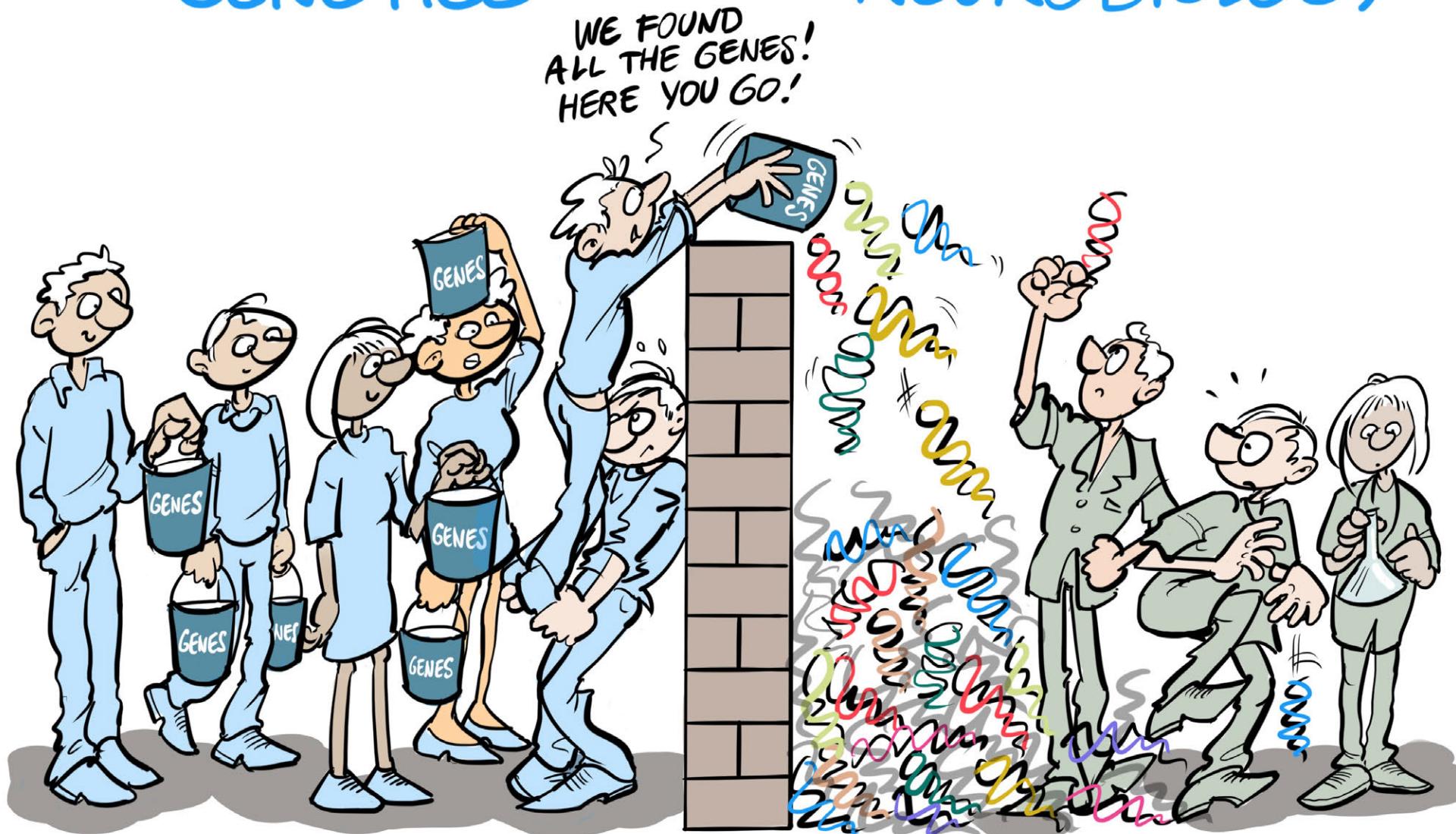
Ahmed Mahfouz

Department of Human Genetics, Leiden University Medical Center  
Pattern Recognition and Bioinformatics, TU Delft

 @ahmedElkoussy  
[mahfouzlab.org](http://mahfouzlab.org)

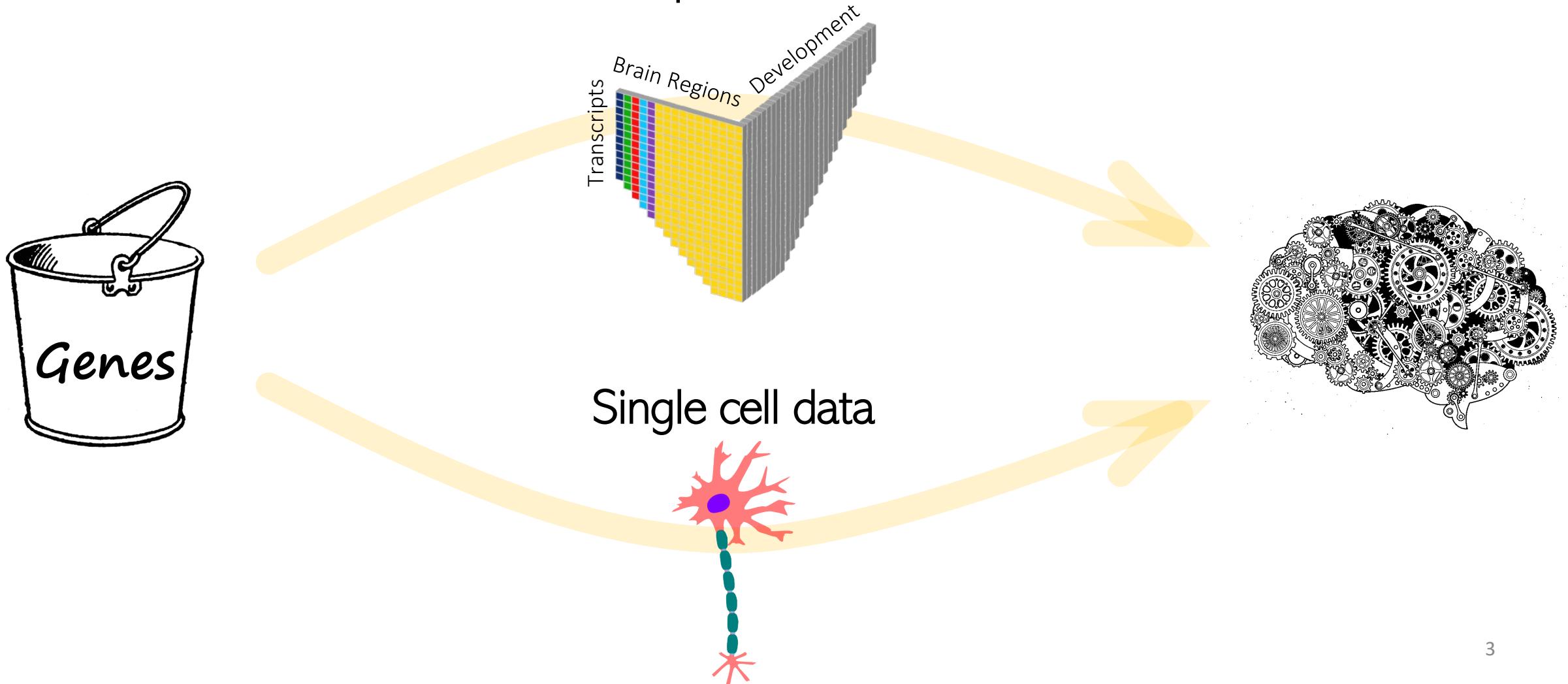
# GENETICS

# NEUROBIOLOGY



# A functional genomics approach

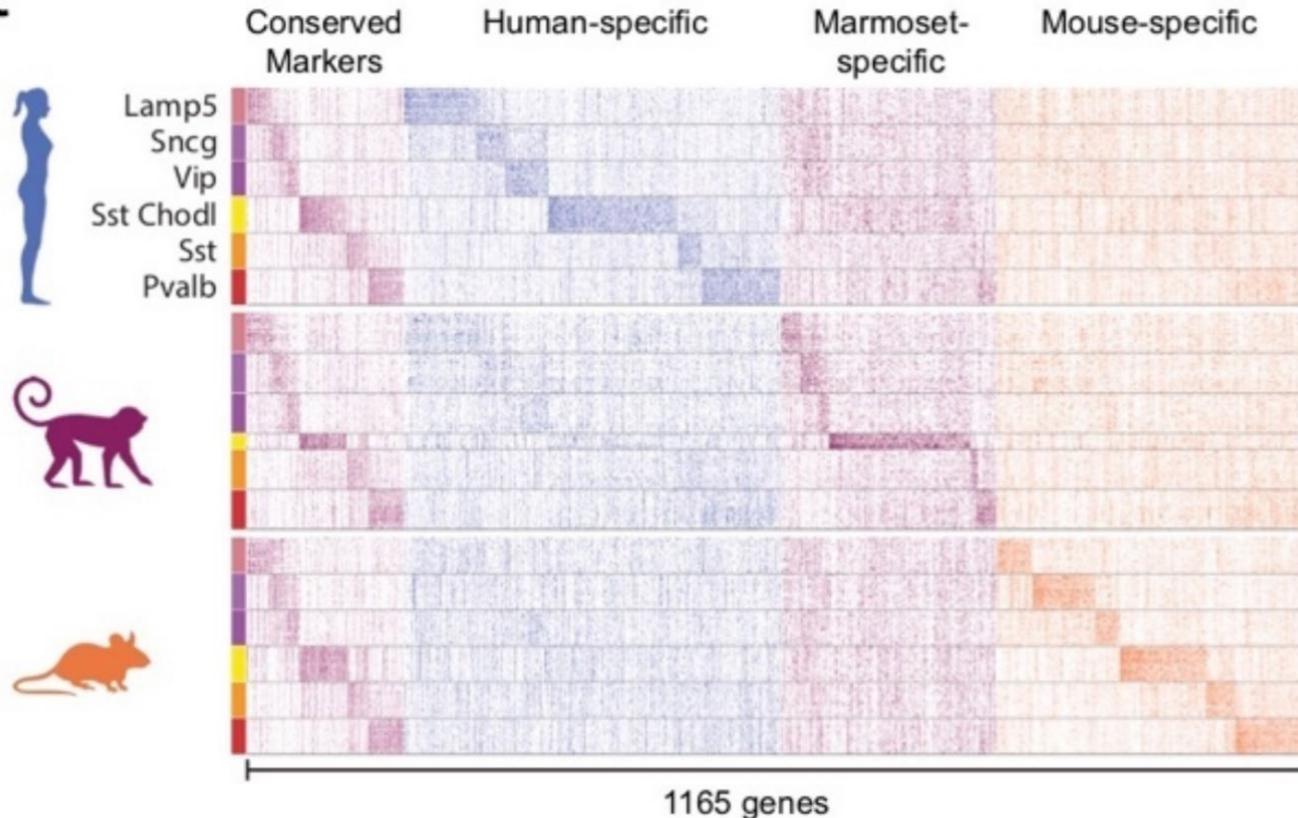
## Transcriptome atlases



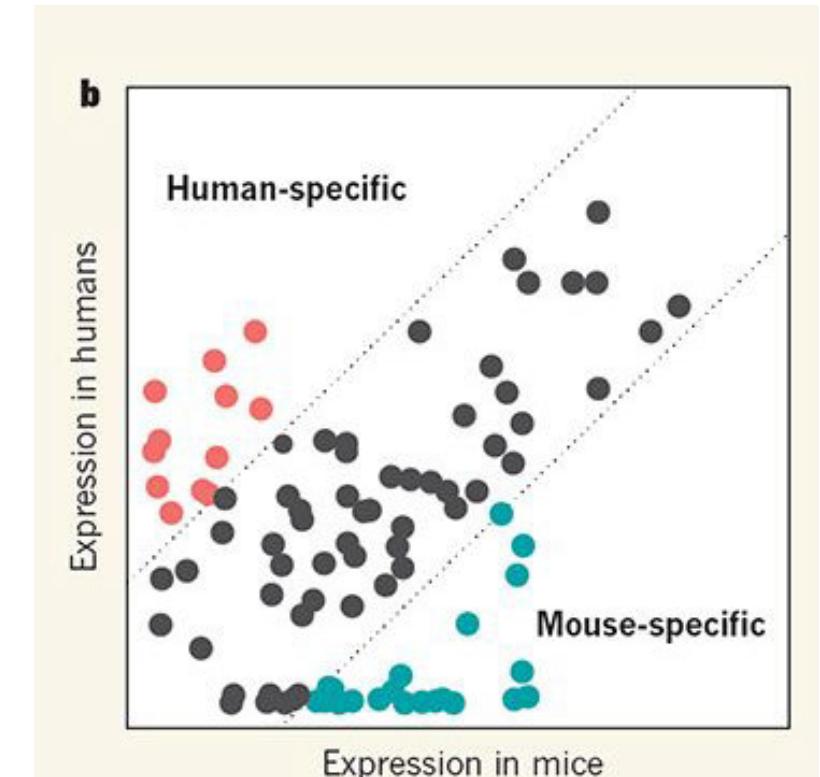
# Comparing cells across studies/species/conditions...

Several marker genes are species-specific

c

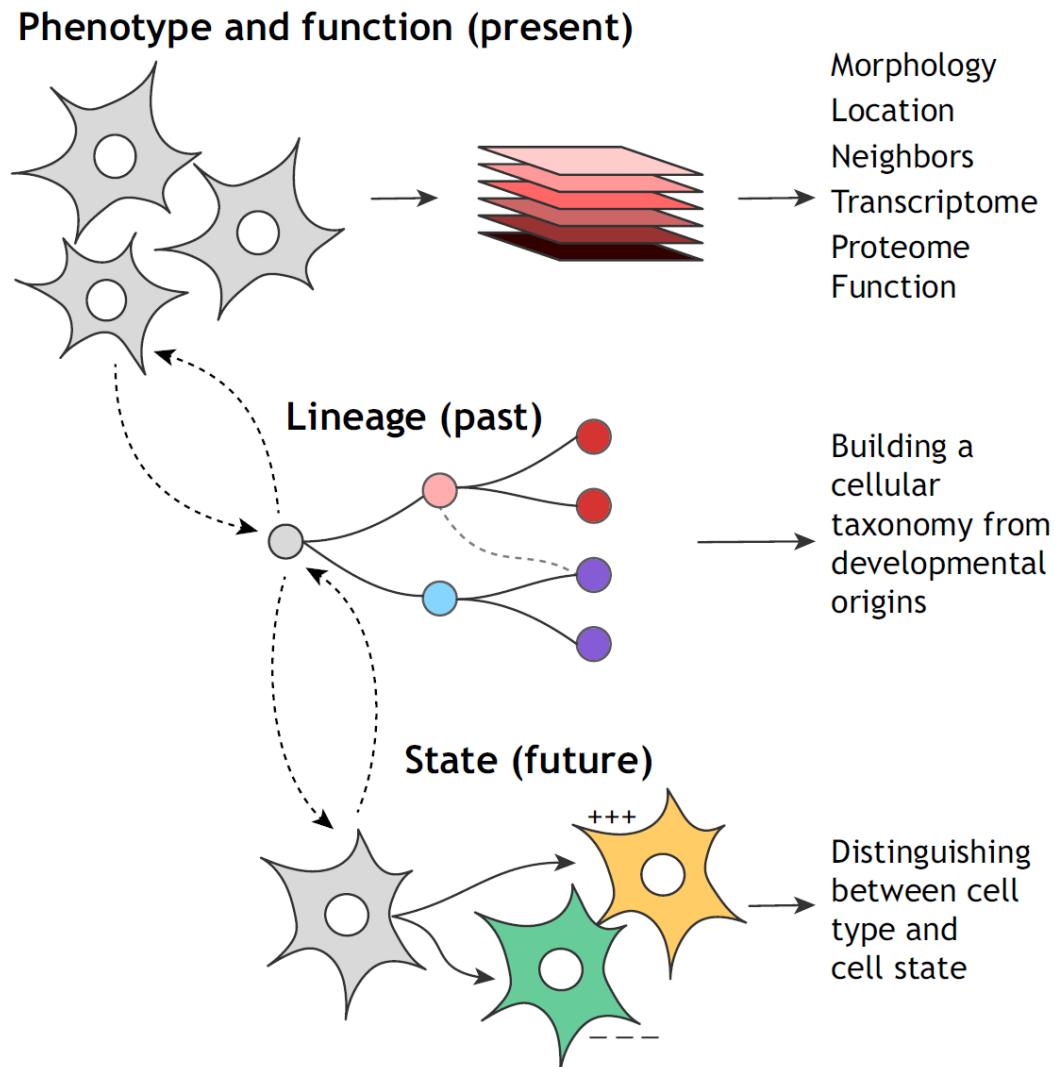


Bakken et al. bioRxiv 2020

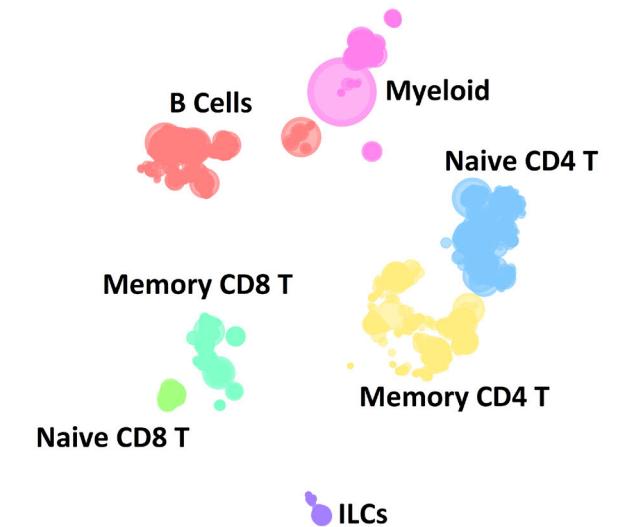
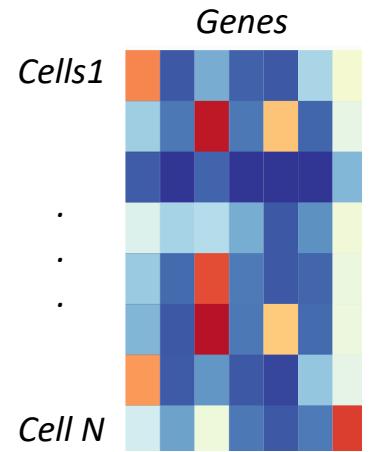


Hodge, Bakken et al. Nature 2019  
Keefe & Nowakowski, Nature 2019

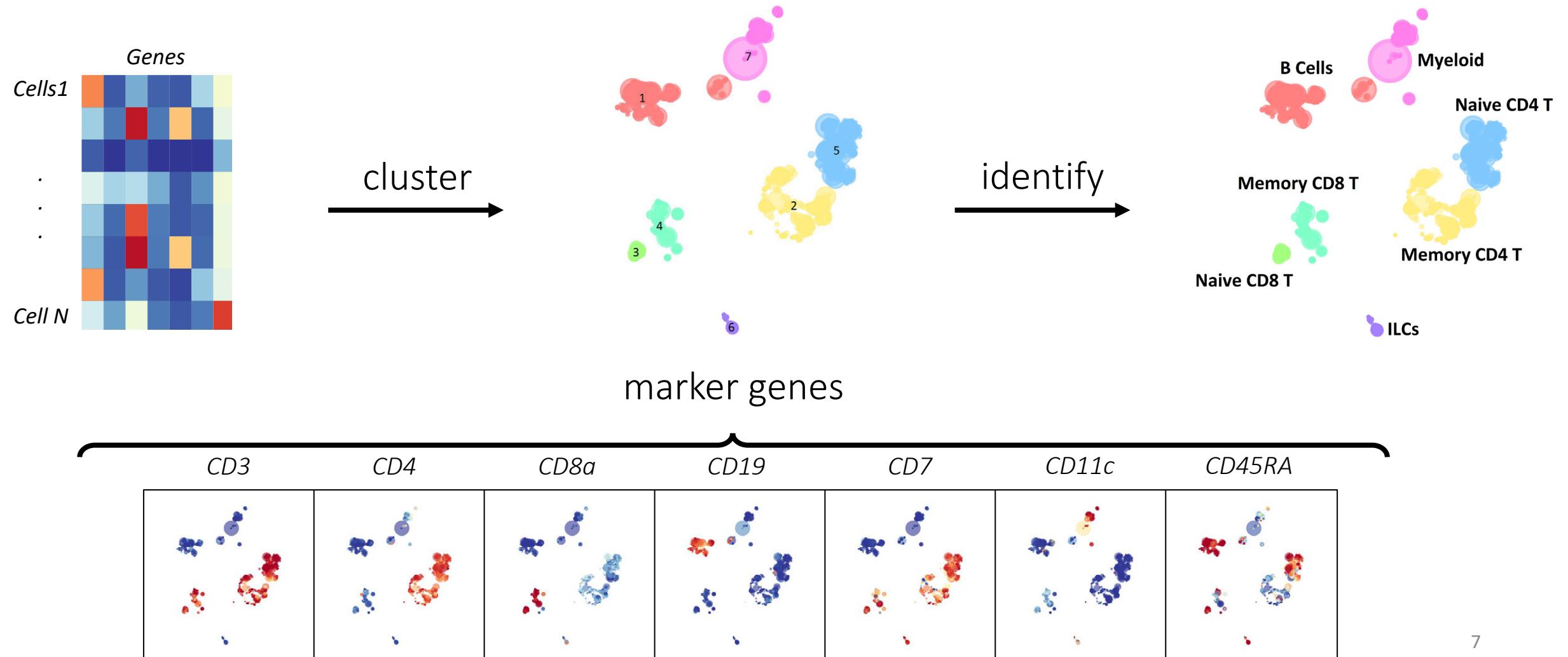
# Cell identity



# How can we identify cell populations?

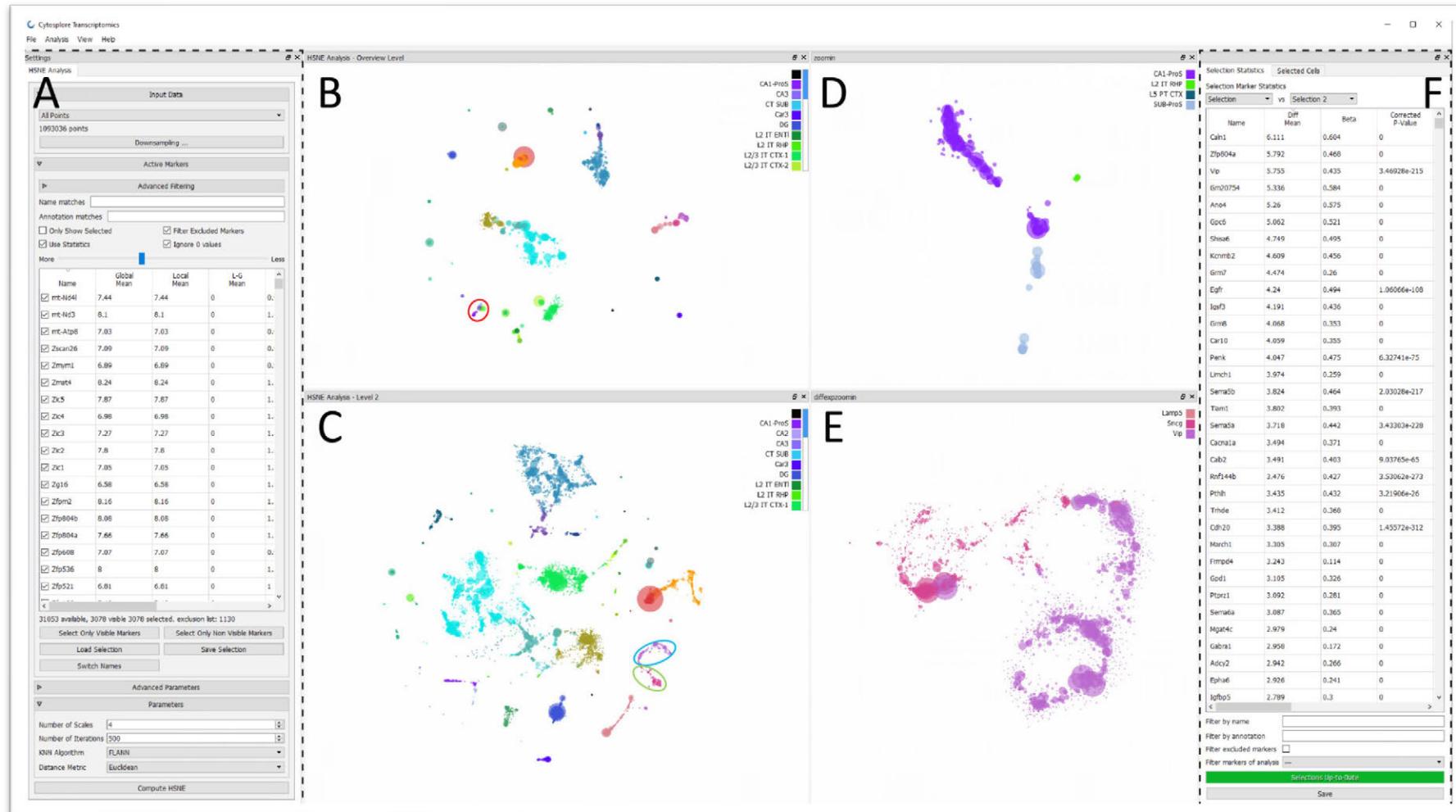


# How can we identify cell populations?





# Cytosplore Transcriptomics

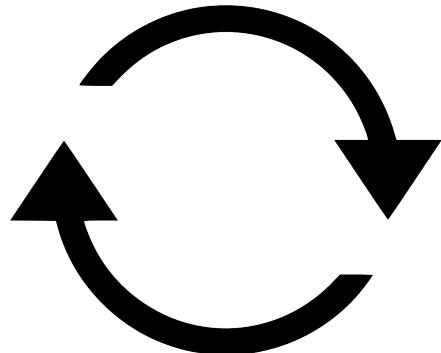


# Unsupervised cell identification is problematic

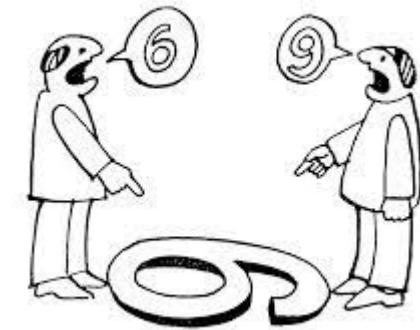
Time consuming



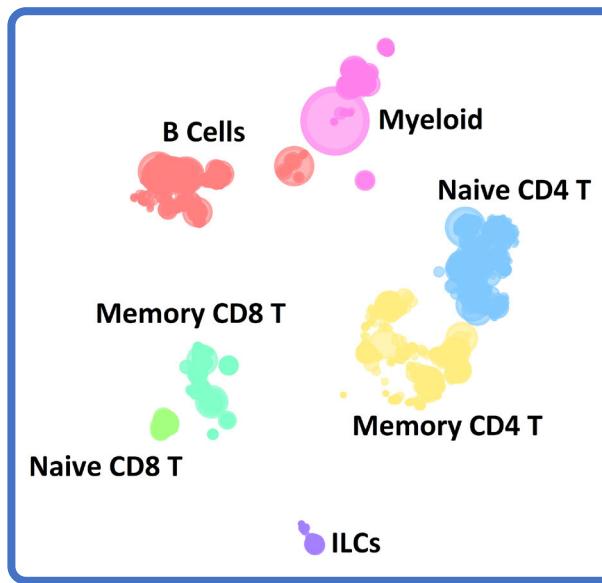
Not reproducible



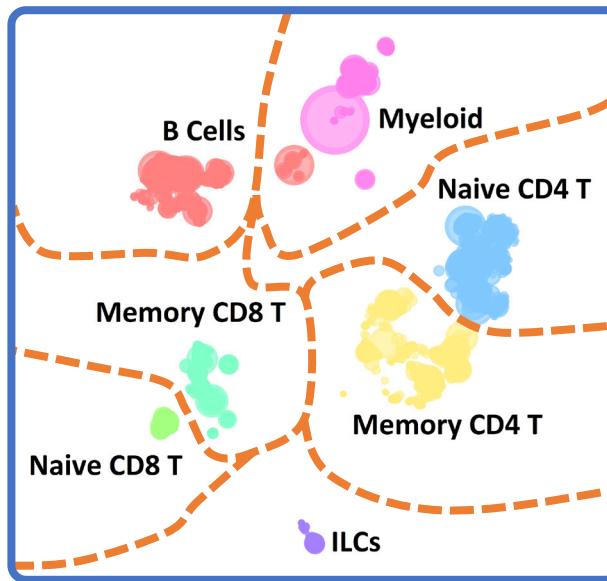
Subjective



# Can we automatically identify cell populations?



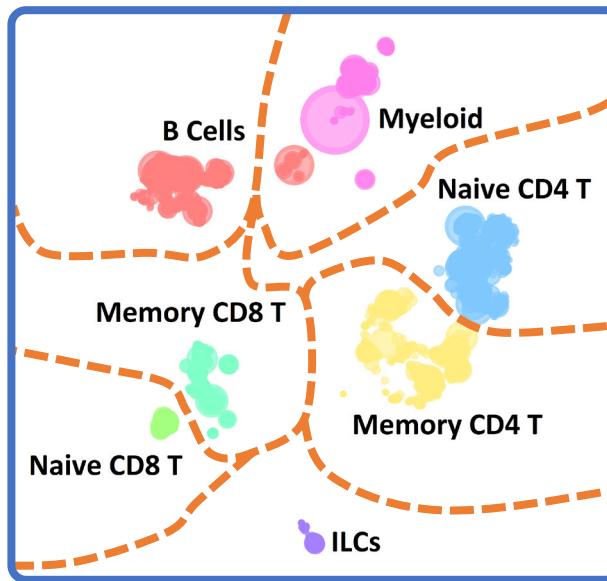
# Can we automatically identify cell populations?



# Can we automatically identify cell populations?

## Clustering

- **Unsupervised** learning
- Discovering structure/relations
- Clusters are defined by a decision boundary

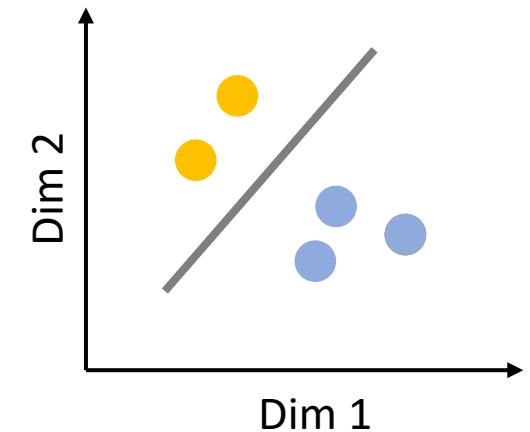


## Classification

- **Supervised** learning
- Prior information available about different groups
- Classifiers find descriptions of decision boundaries

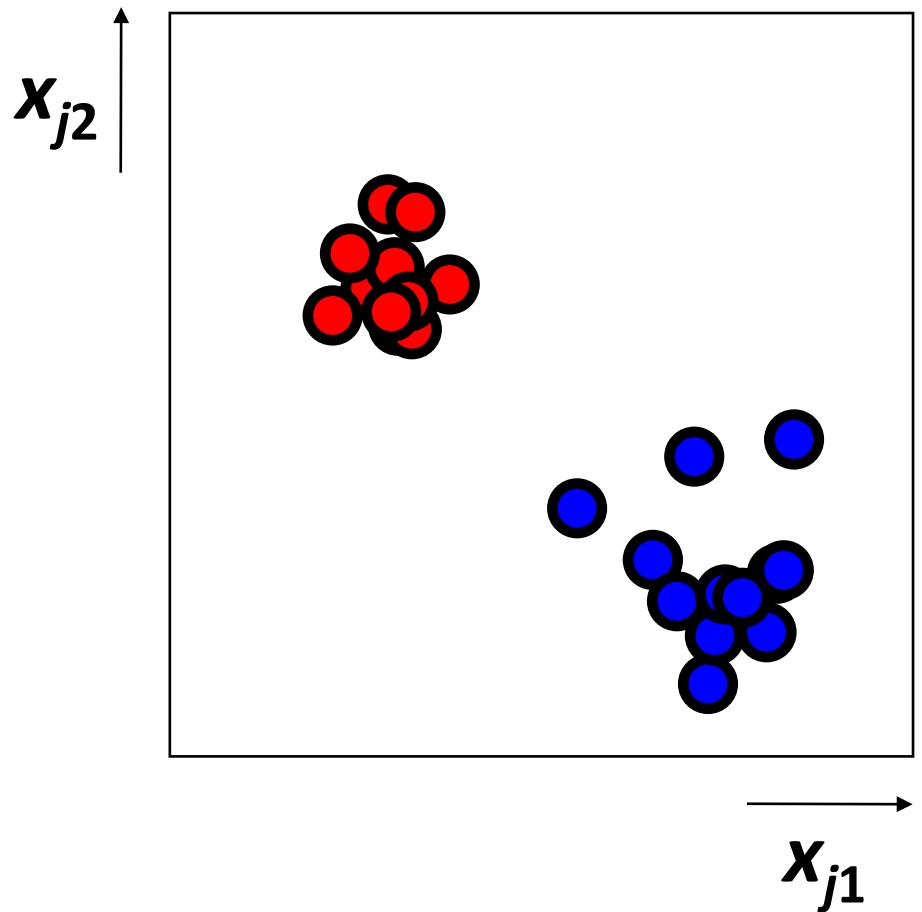
# Classification

	Features (genes)												Labels (cell types)	
Cell 1													T cell	
Cell 2													T cell	
Cell N													B cell	
													B cell	
Unknown Cell													T cell	



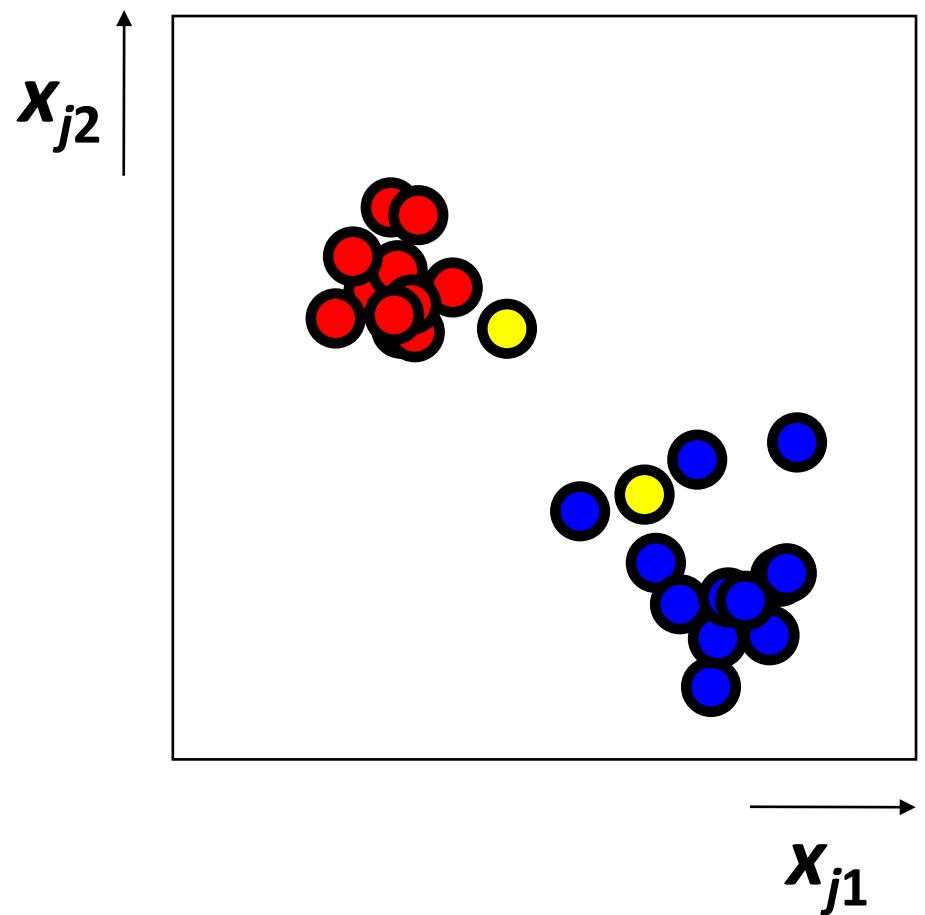
# Classifier training

- Dataset: for  $j^{\text{th}}$  cell:
  - gene expressions  $\mathbf{x}_j$
  - class label:  $y_j \in \{1=\text{T}, -1=\text{B}\}$
- Classifier:  $\hat{y}_j = W(\mathbf{x}_j)$
- Errors:  $E = \text{sum}(E_j)$      $E_j = \begin{cases} 1 & \text{if } \hat{y}_j \neq y_j \\ 0 & \text{if } \hat{y}_j = y_j \end{cases}$
- Place decision boundary (i.e. change  $W$ ) s.t.  $E$  is minimal



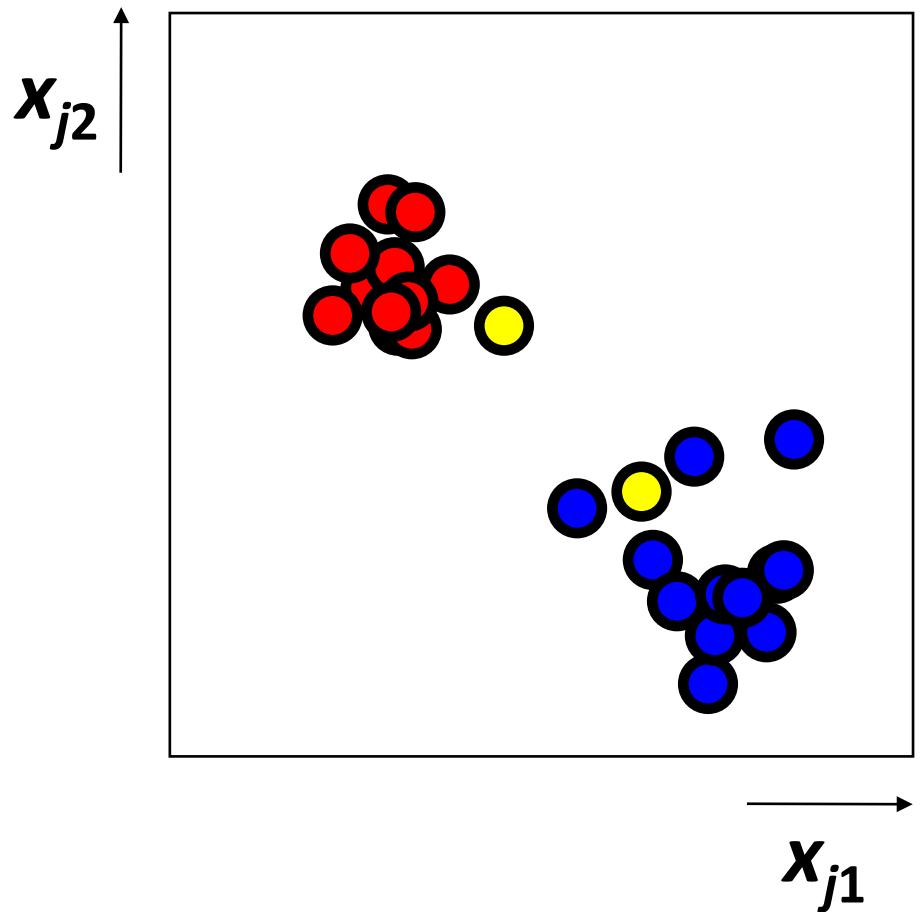
# Instance Based Learning (Lazy Classification)

- Example: **Nearest neighbor (k-NN)**
  - Keep the whole training dataset
  - A query example (vector) comes
  - Find closest example(s)
  - Predict
- *No actual training*



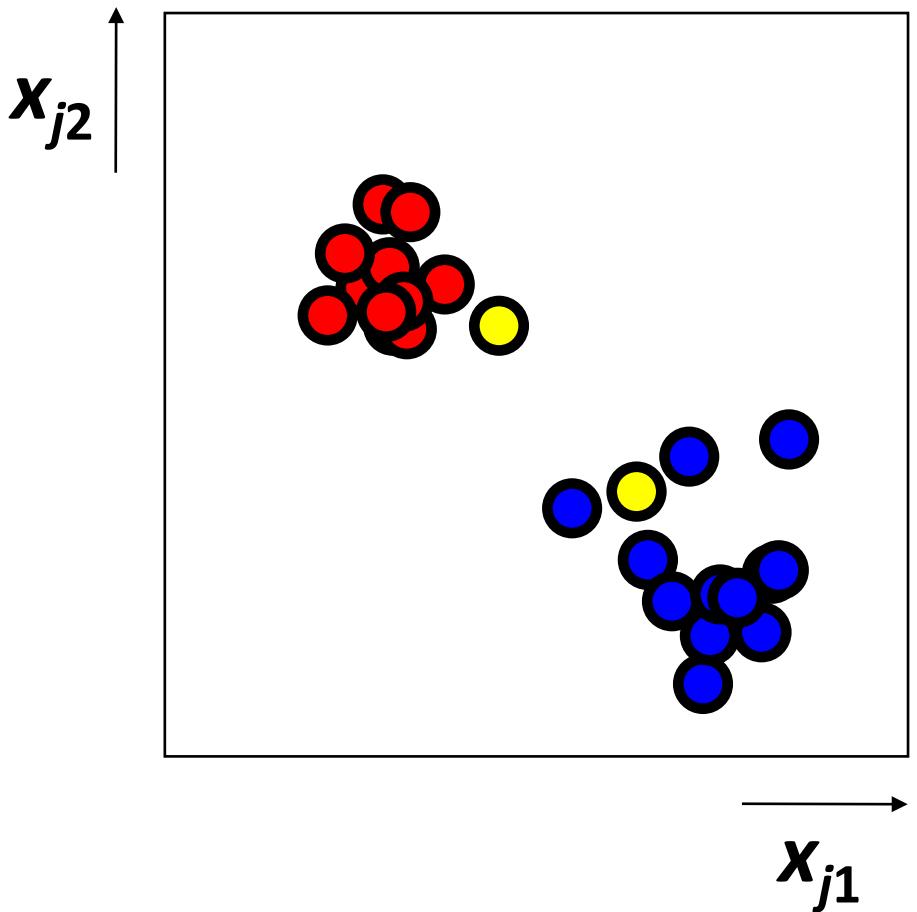
# Nearest Neighbor (k-NN)

- To make Nearest Neighbor work we need 4 things:
  - 1) Distance metric:
  - 2) How many neighbors to look at?
  - 3) Weighting function (optional)
  - 4) How to fit with the local points?



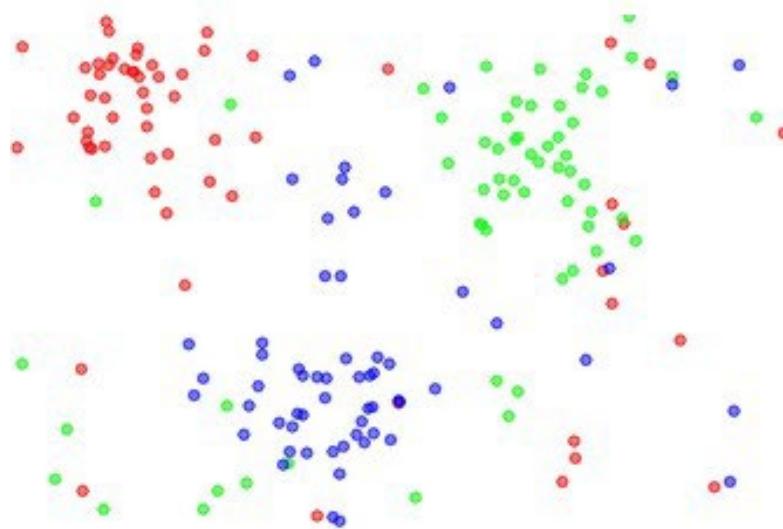
# Nearest Neighbor (k-NN)

- Distance metric:
  - Euclidean
- How many neighbors to look at?
  - $k$
- Weighting function (optional):
  - Unused
- How to fit with the local points?
  - Predict the average output among  $k$  nearest neighbors

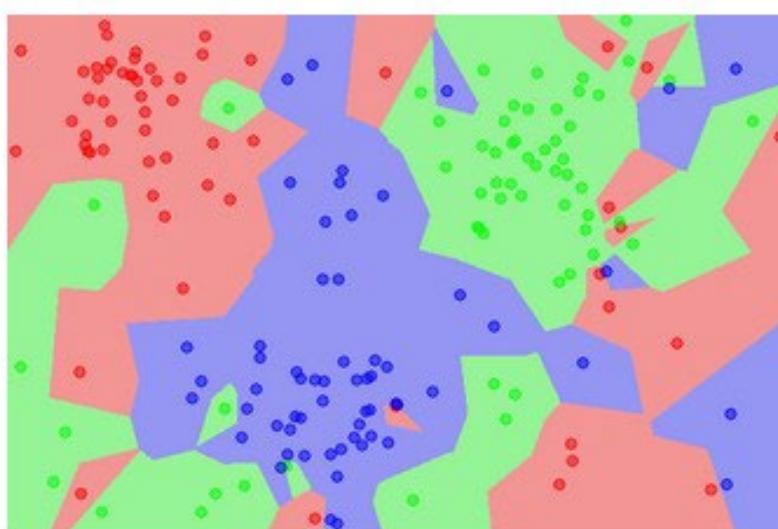


# Effect of $k$

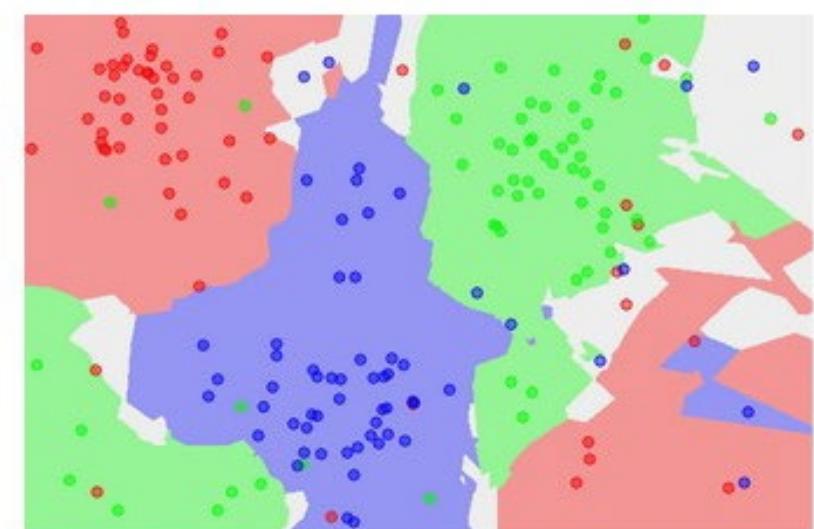
the data



NN classifier

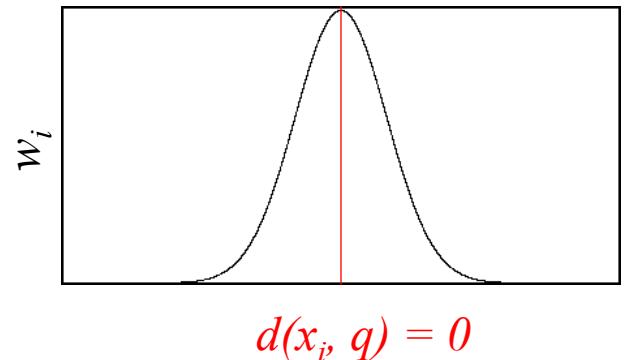


5-NN classifier

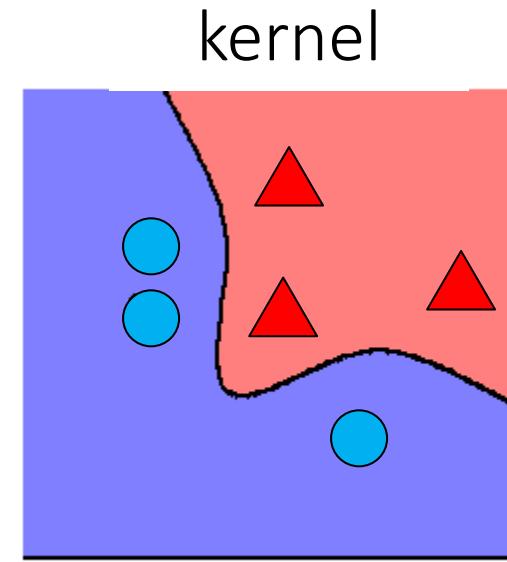
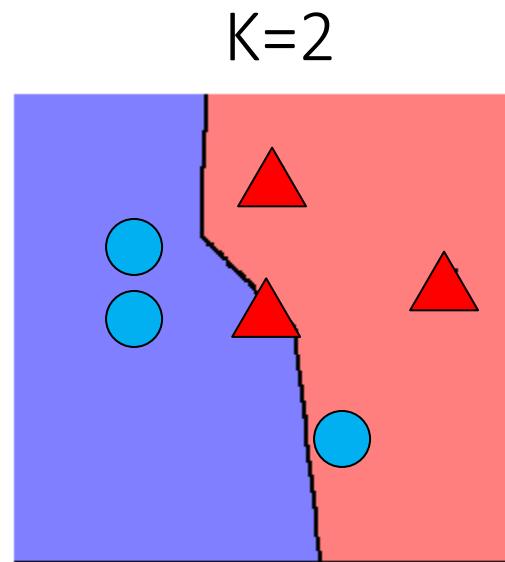
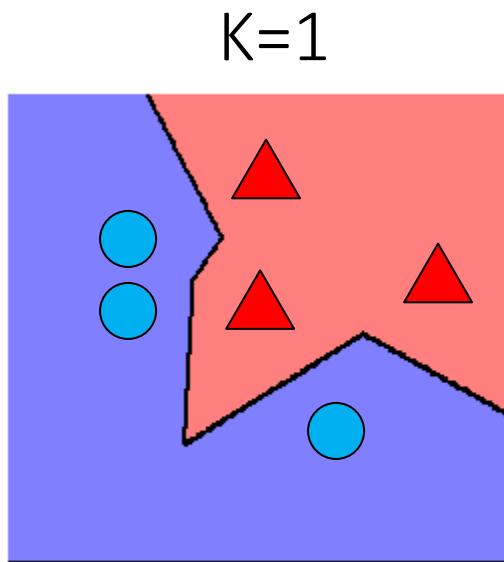


# Weighted Nearest Neighbor (kernel regression)

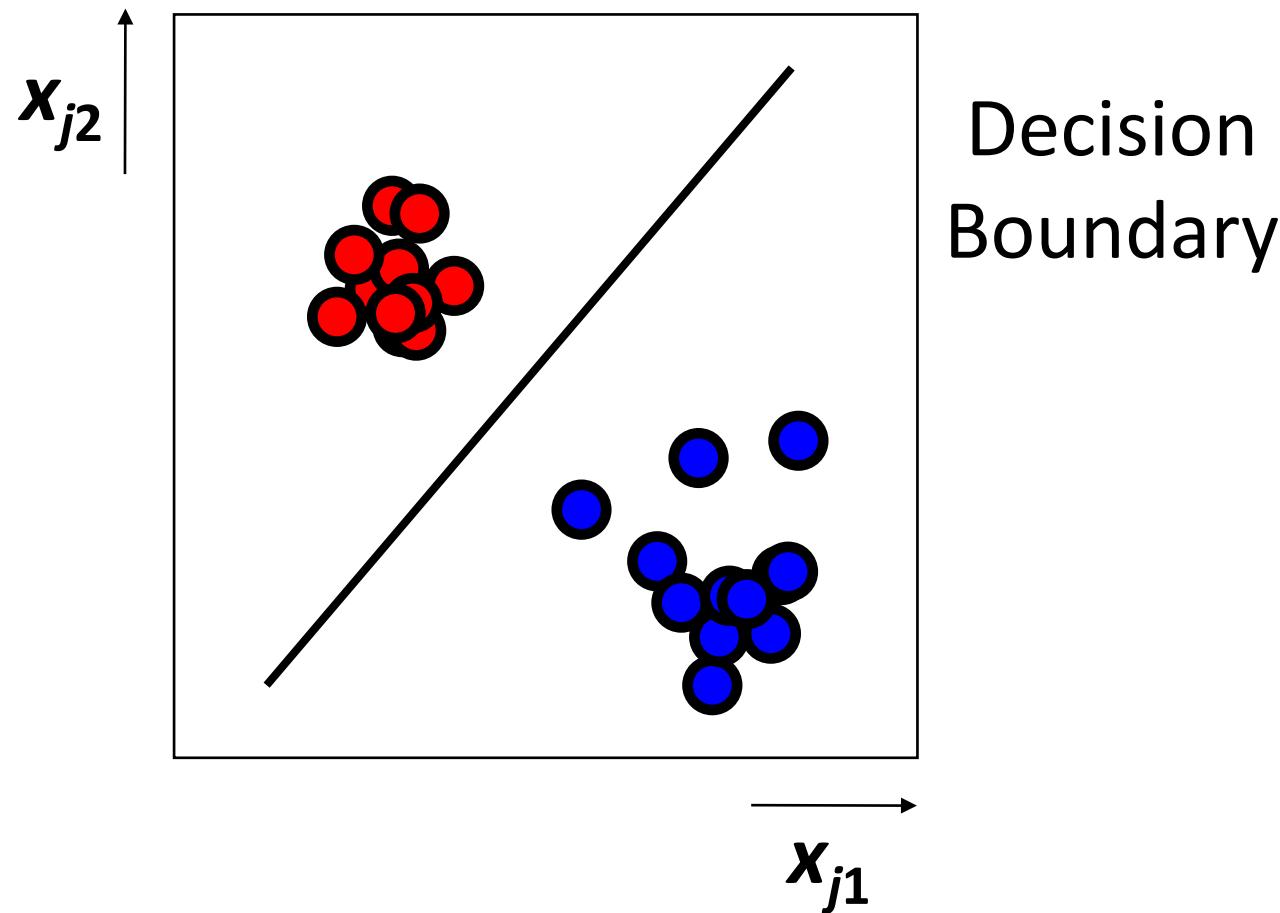
- Distance metric:
  - Euclidean
- How many neighbors to look at?
  - All of them (!)
- Weighting function:
  - $w_i = \exp\left(-\frac{d(x_i, q)^2}{K_w}\right)$
  - Nearby points to query  $q$  are weighted more strongly.  $K_w$ : kernel width.
- How to fit with the local points?
  - Predict weighted average:  $\frac{\sum_i w_i y_i}{\sum_i w_i}$



# Comparison: K=1, K=2, kernel

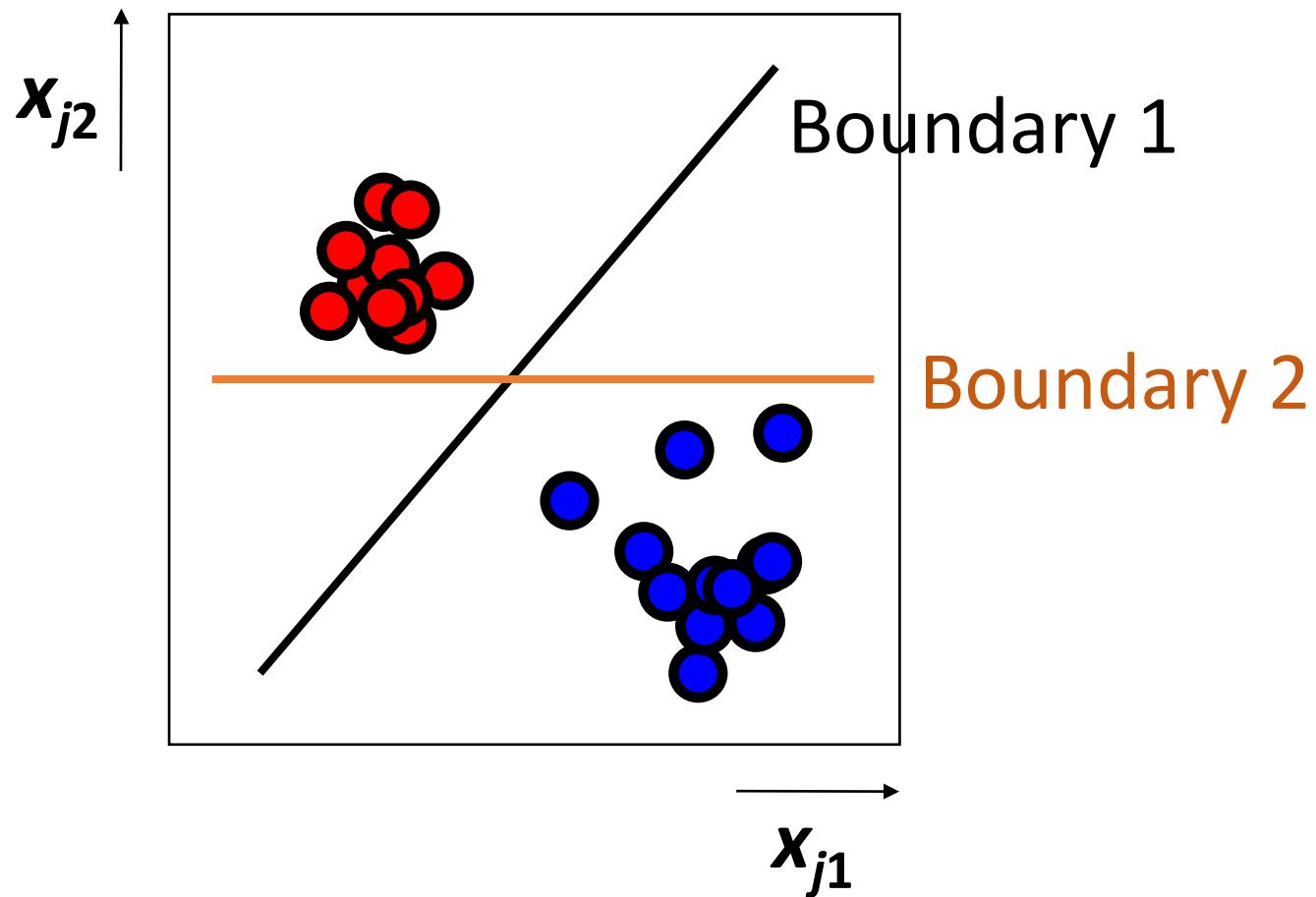


# Support Vector Machine (SVM)



# Support Vector Machine (SVM)

Which boundary is better?

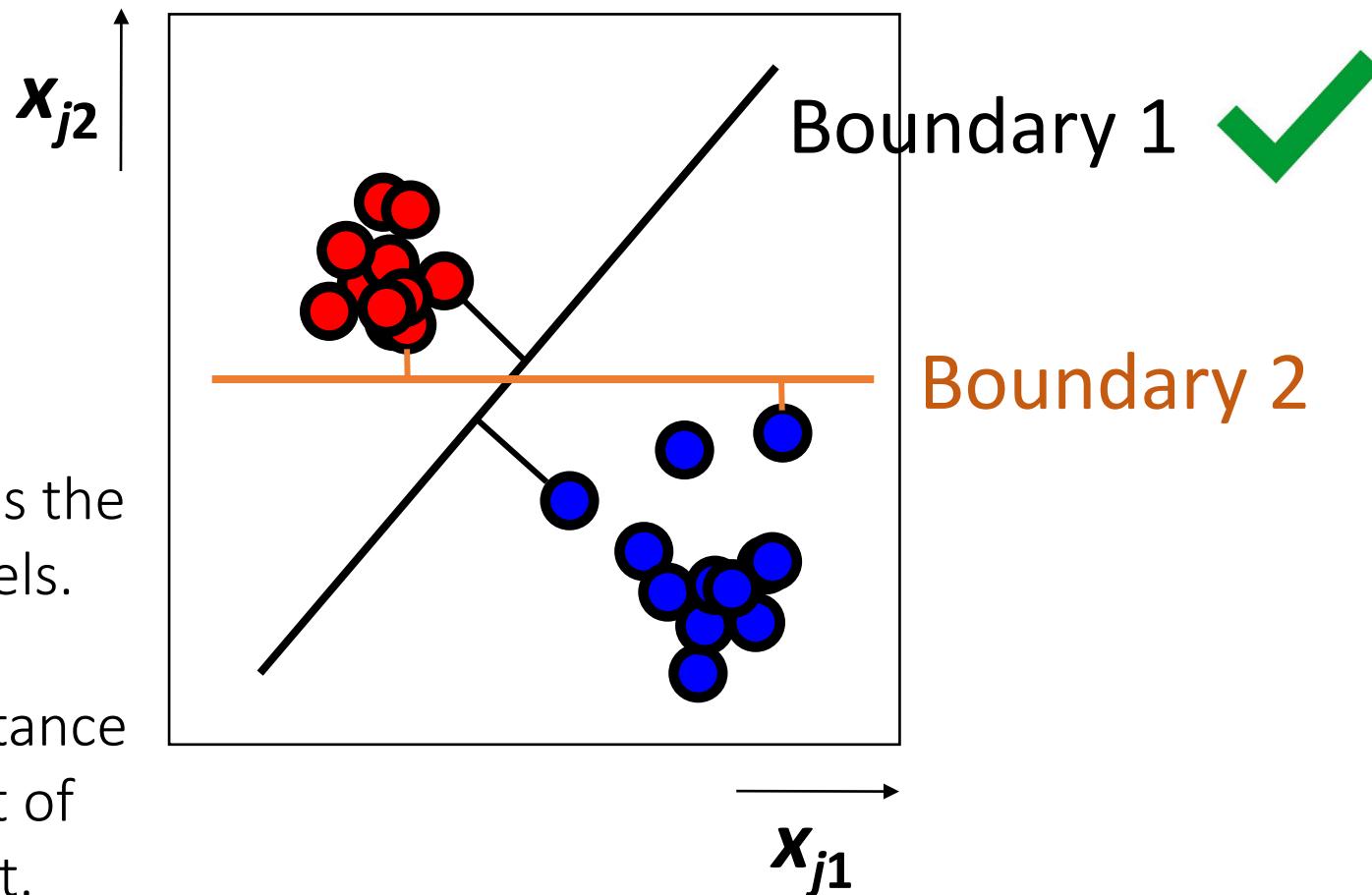


# Support Vector Machine (SVM)

Which boundary is better?

The one that maximizes the margins from both labels.

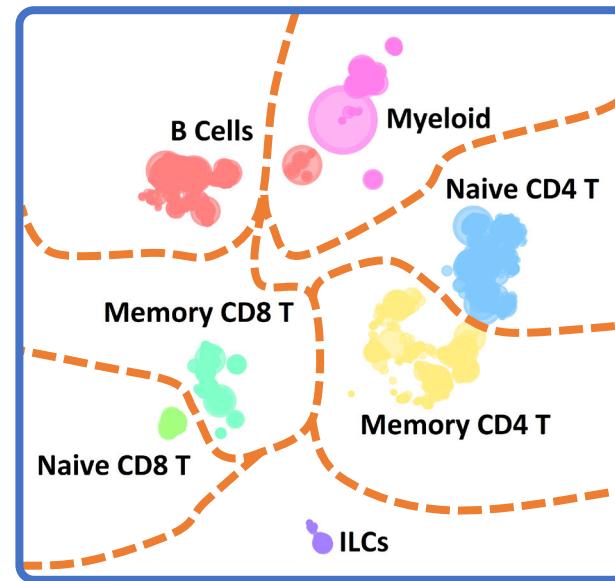
i.e. The one whose distance to the nearest element of each label is the largest.



# Can we automatically identify cell populations?

*Training data*

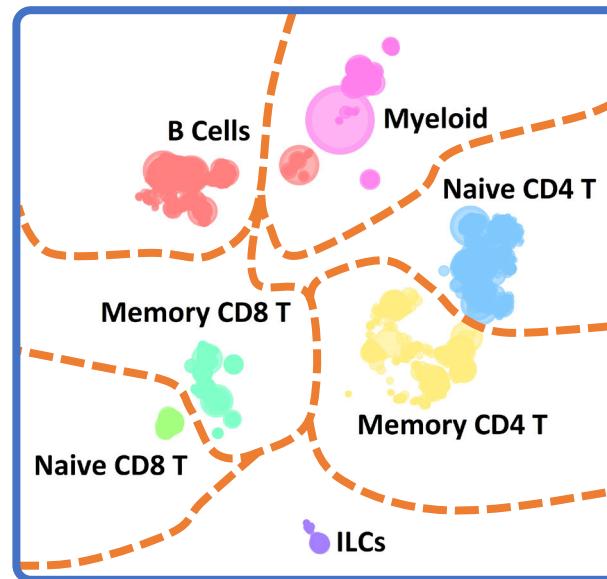
Annotated Cells  
(e.g. atlas)



# Can we automatically identify cell populations?

*Training data*

Annotated Cells  
(e.g. atlas)

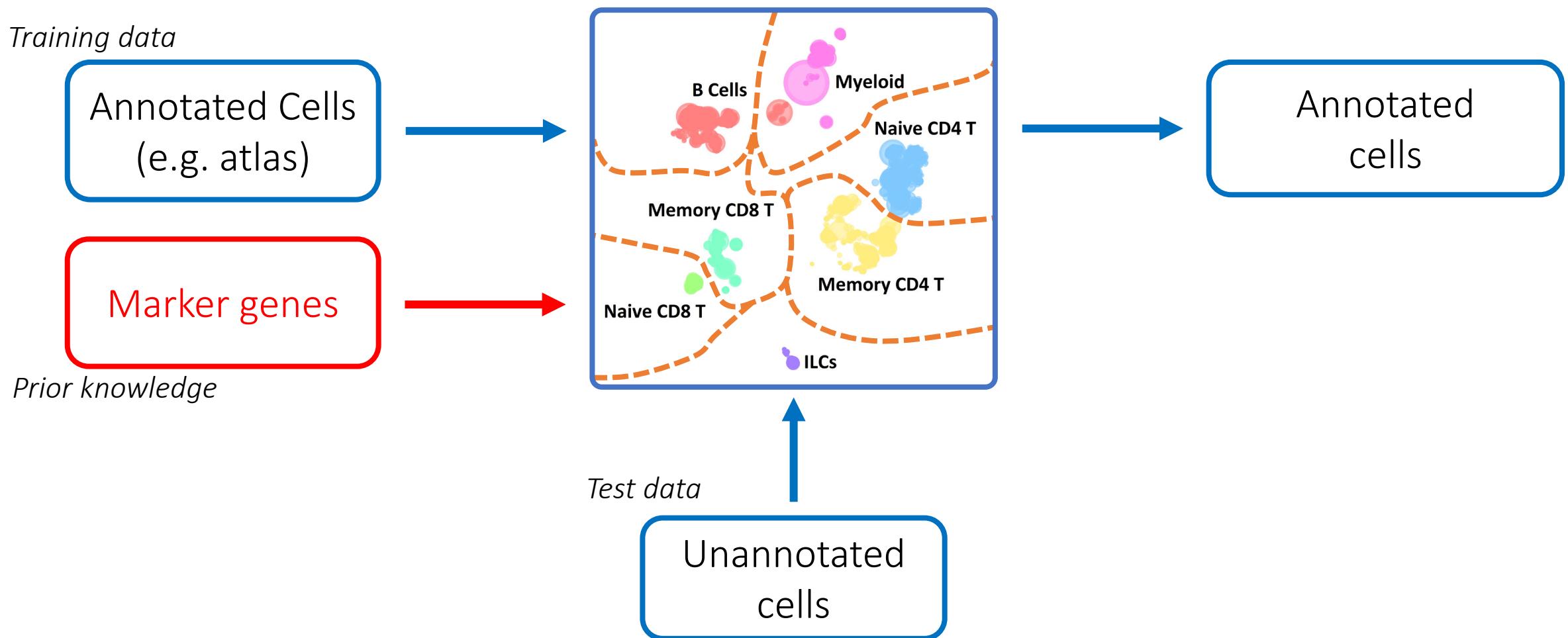


*Prior knowledge*

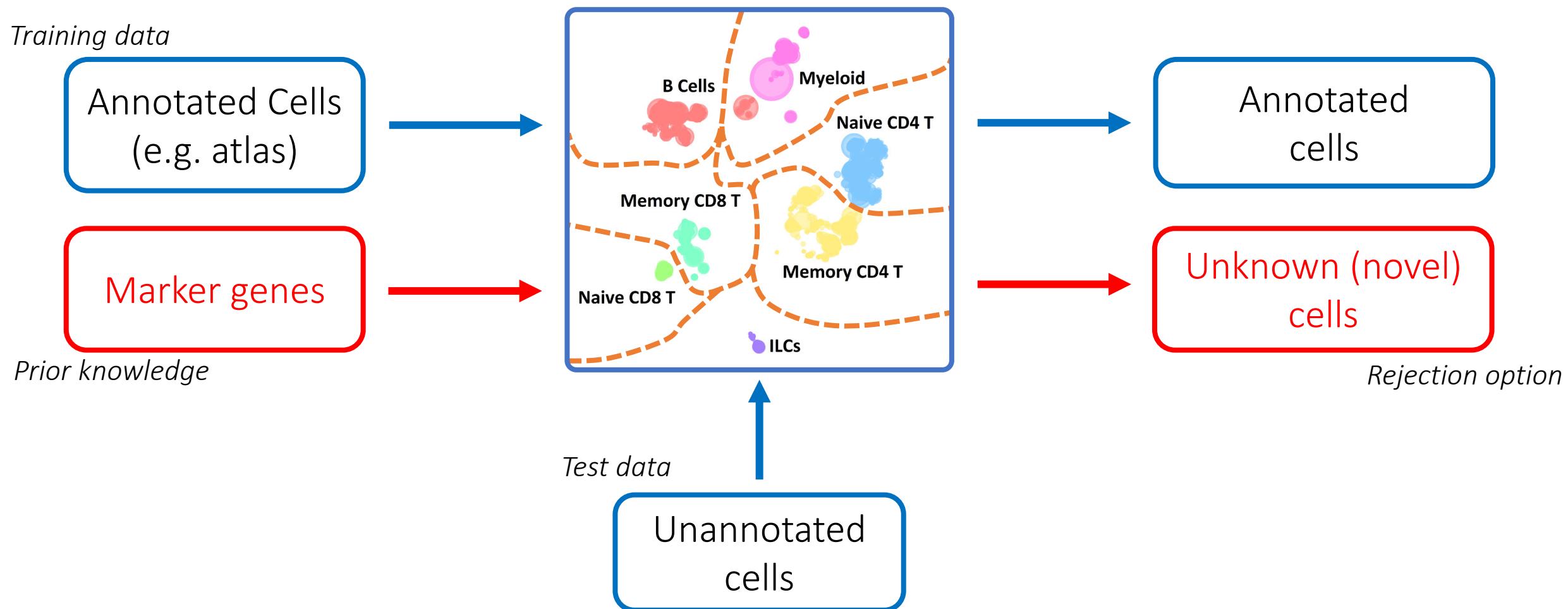
Marker genes



# Can we automatically identify cell populations?



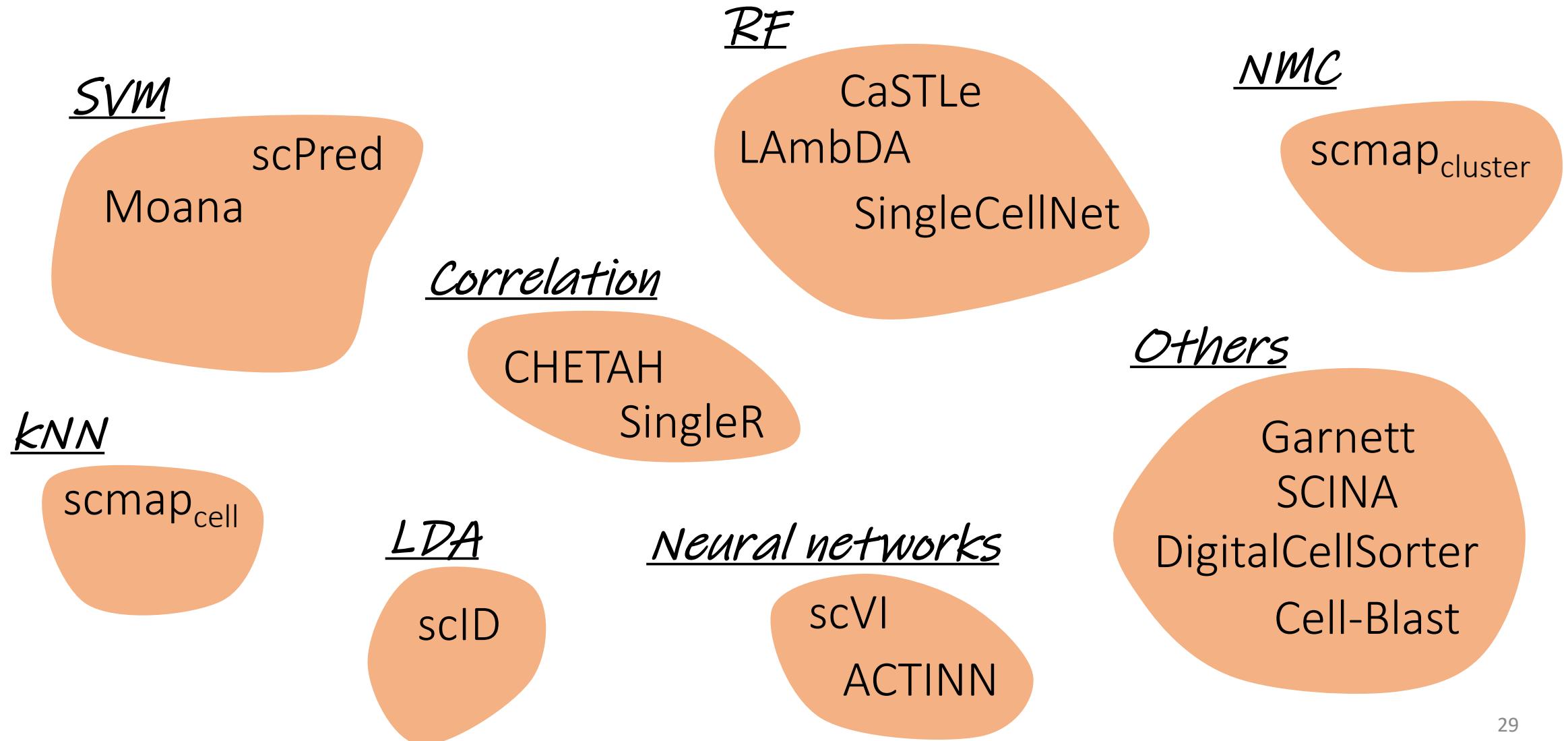
# Can we automatically identify cell populations?



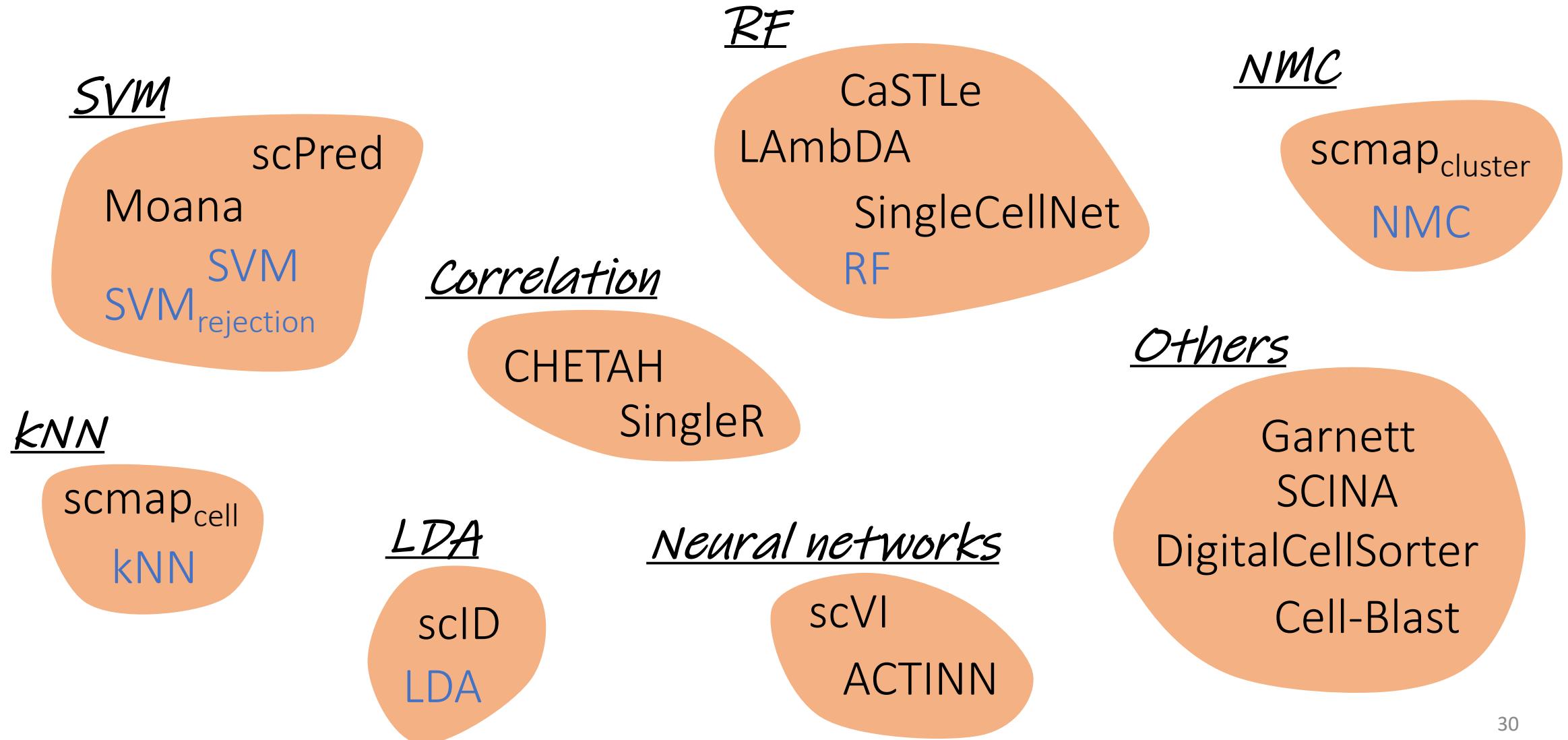
# 16 existing classifiers (April 2019)

scPred	CaSTLe	scmap <sub>cluster</sub>
Moana	LAmbDA	
	SingleCellNet	
scmap <sub>cell</sub>	CHETAH	Garnett
	SingleR	SCINA
		DigitalCellSorter
scID	scVI	Cell-Blast
	ACTINN	

# 16 existing classifiers (April 2019)

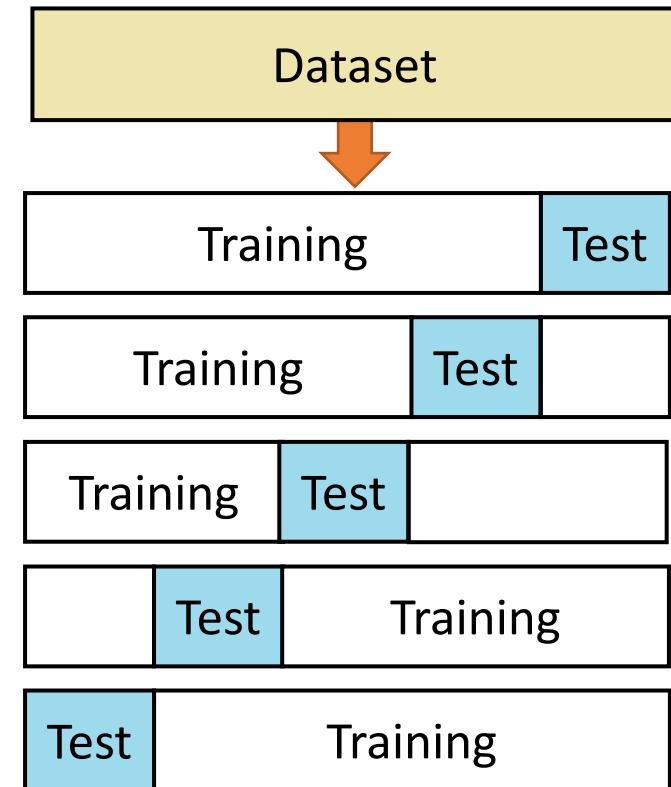


# 16 existing + 6 off-the-shelf classifiers



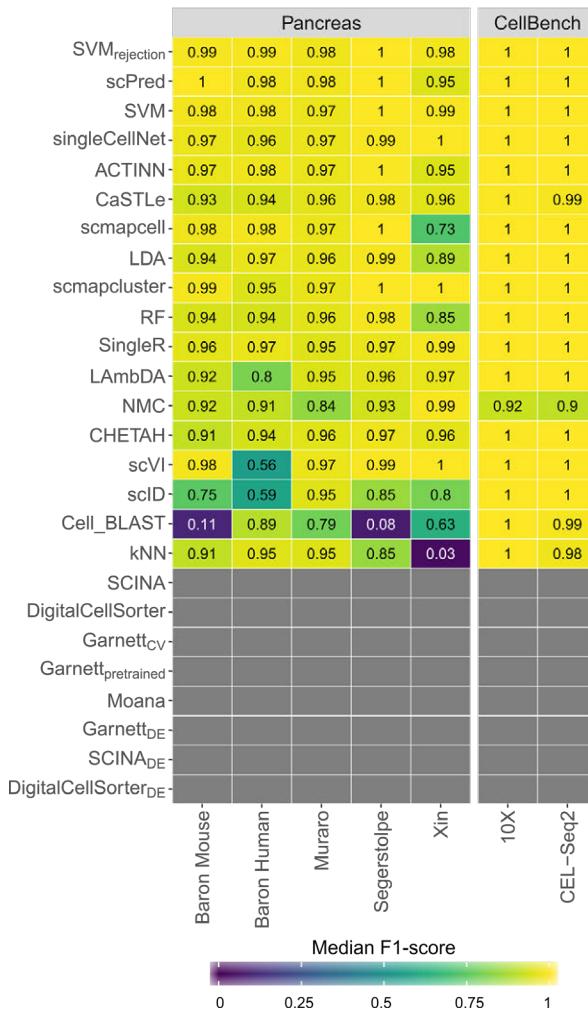
# Experiment 1: intra-dataset evaluation

- Stratified 5-fold cross validation
- Performance evaluation
  - Median F1-score:  $F1 = 2 \frac{precision \cdot recall}{precision + recall}$
  - % unlabelled cells



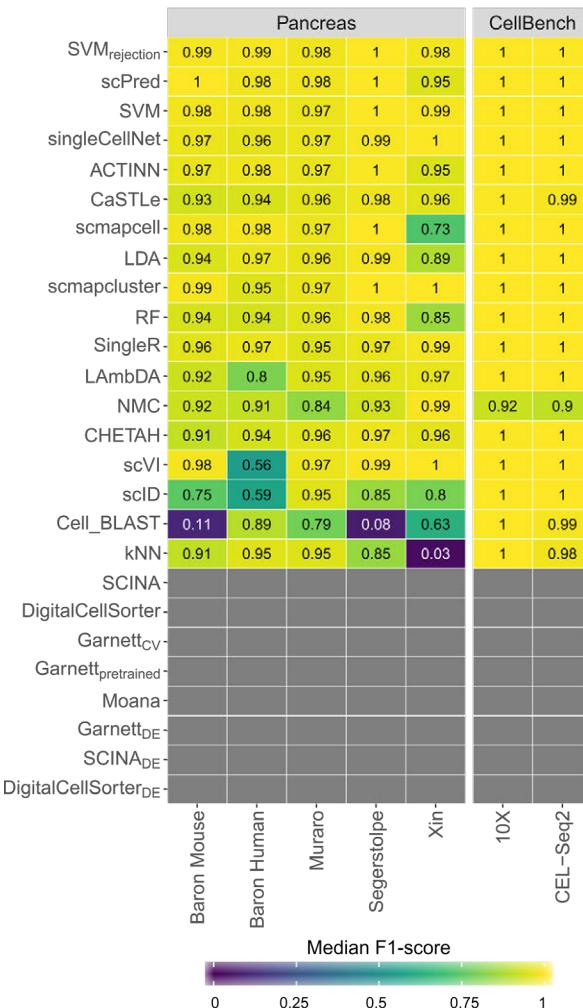
# Most classifiers work well

Median F1-score

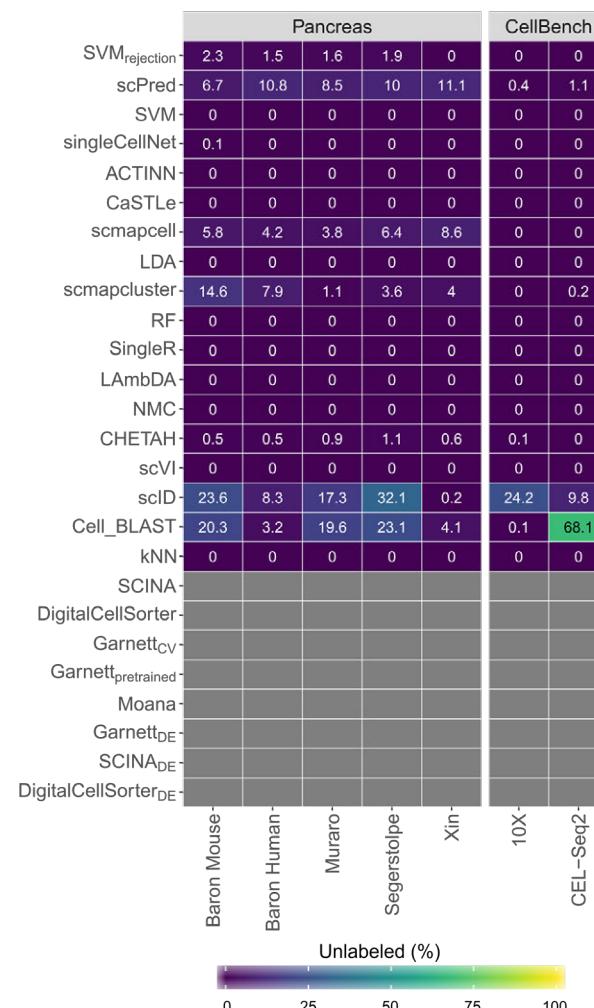


# Most classifiers work well

Median F1-score

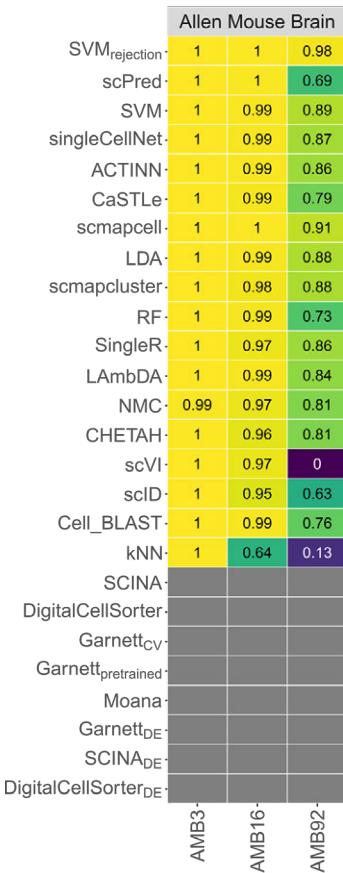


% Unlabeled



# Performance drops with deeper annotation

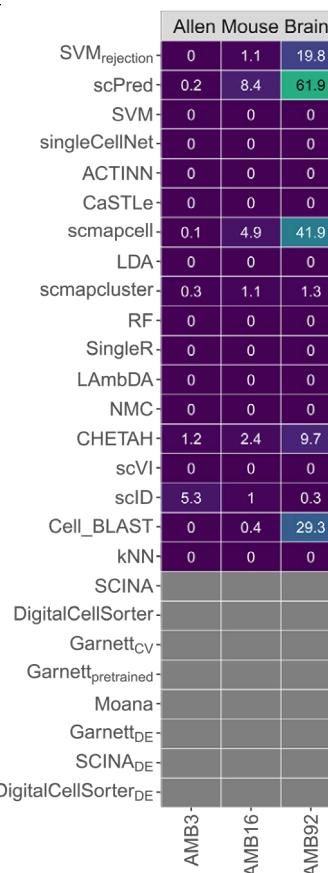
Median F1-score



Median F1-score



% Unlabeled

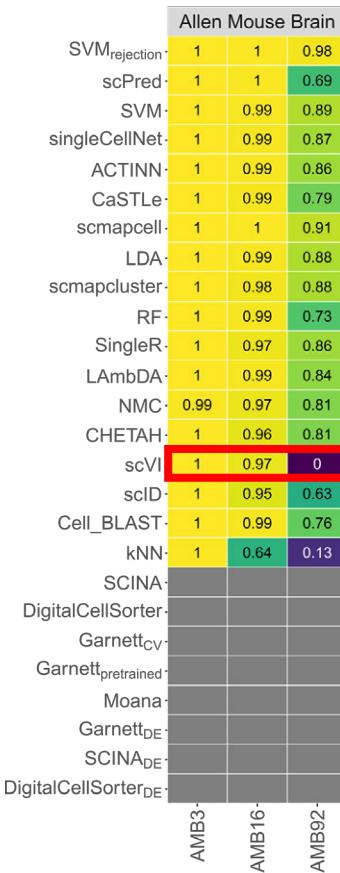


Unlabeled (%)

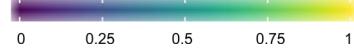


# Performance drops with deeper annotation

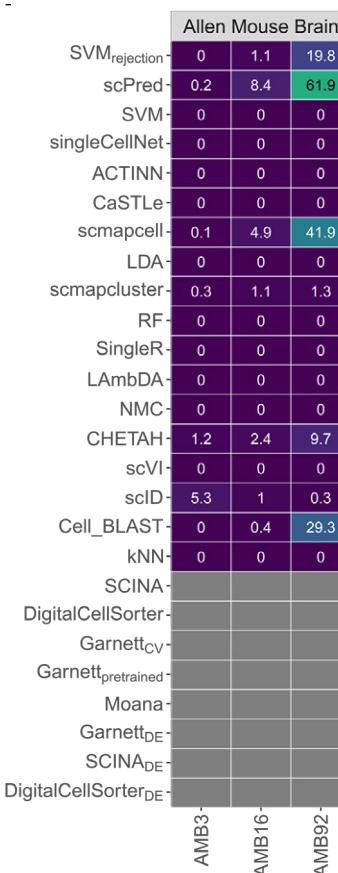
Median F1-score



Median F1-score



% Unlabeled

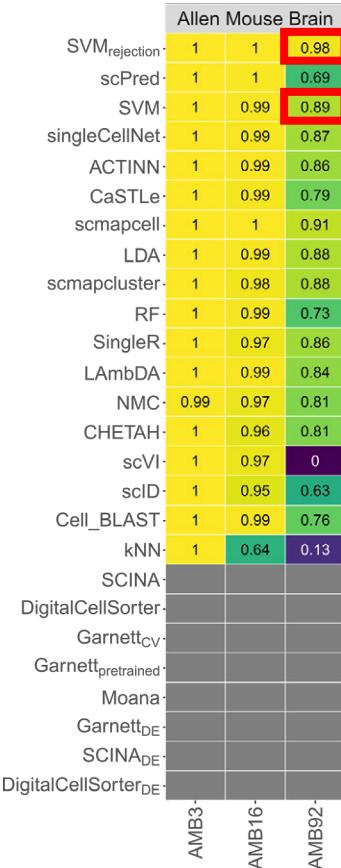


Unlabeled (%)

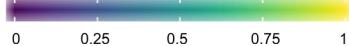


# Trade-off between high performance and rejecting cells

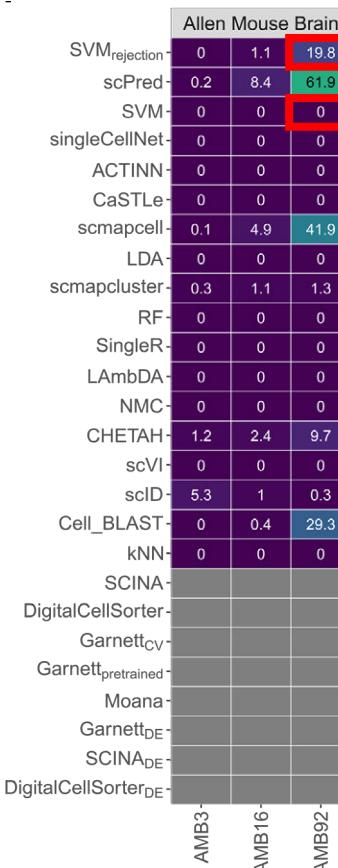
Median F1-score



Median F1-score



% Unlabeled

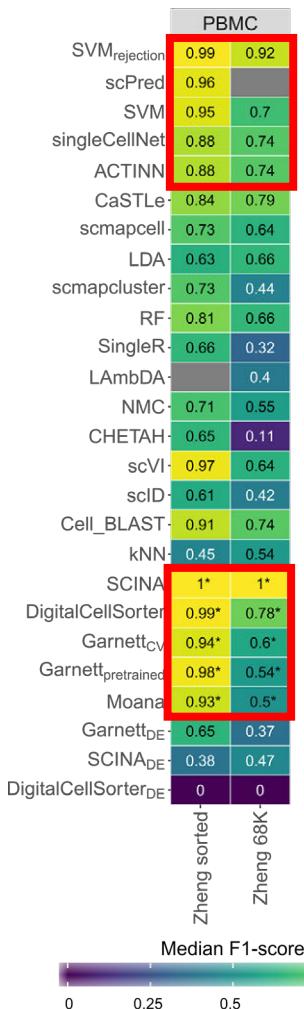


Unlabeled (%)

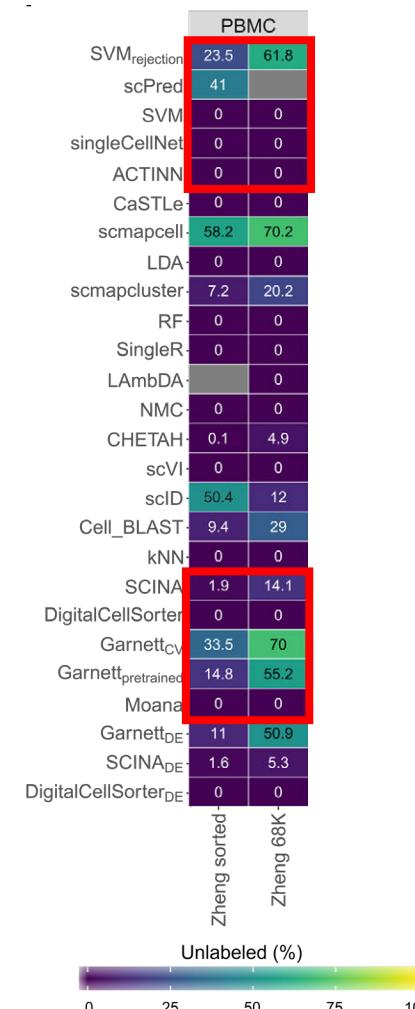


# Prior knowledge is not beneficial

Median F1-score



% Unlabeled



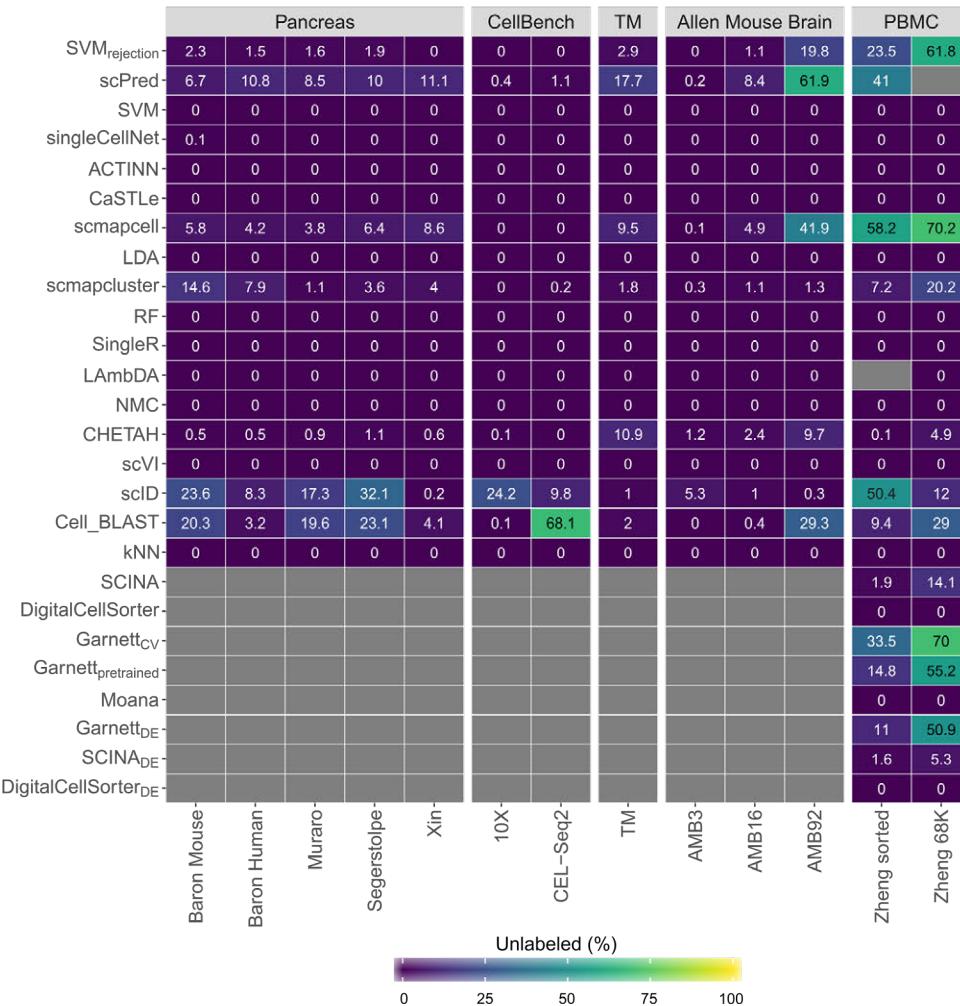
Lower  
number of  
classes!

# Off-the-shelf SVM outperforms dedicated single cell classifiers

Median F1-score

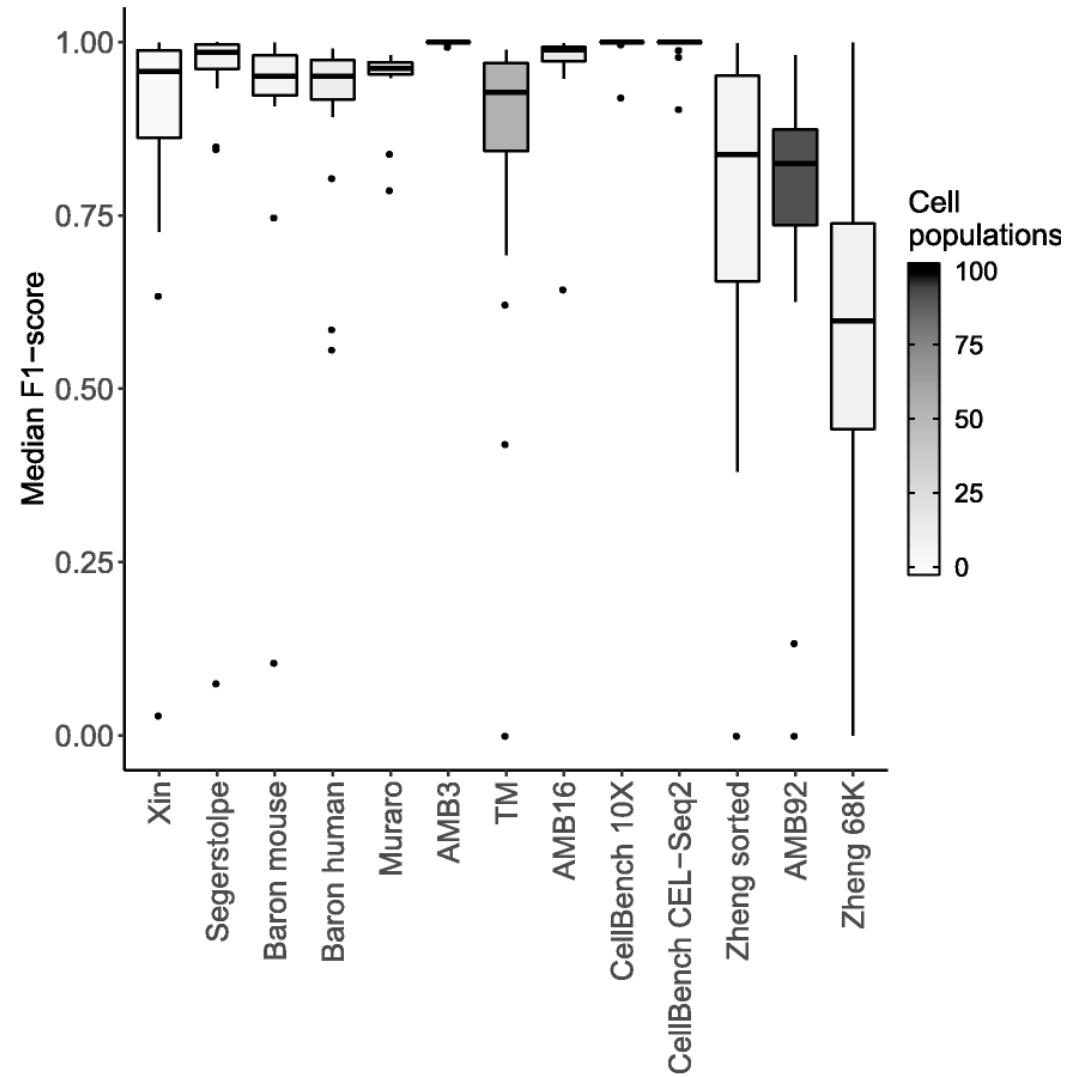
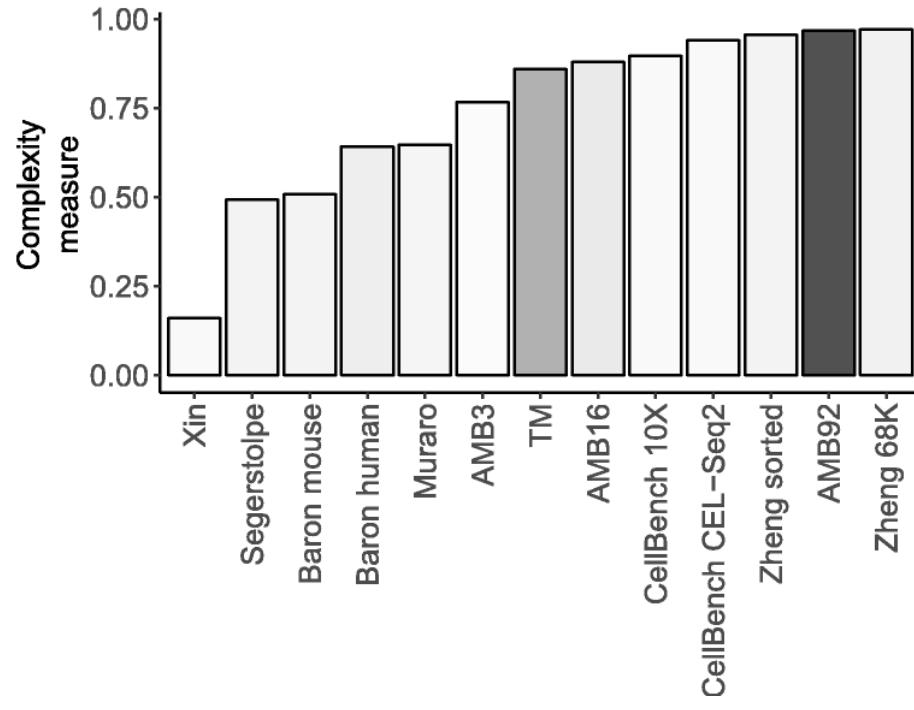


% Unlabeled



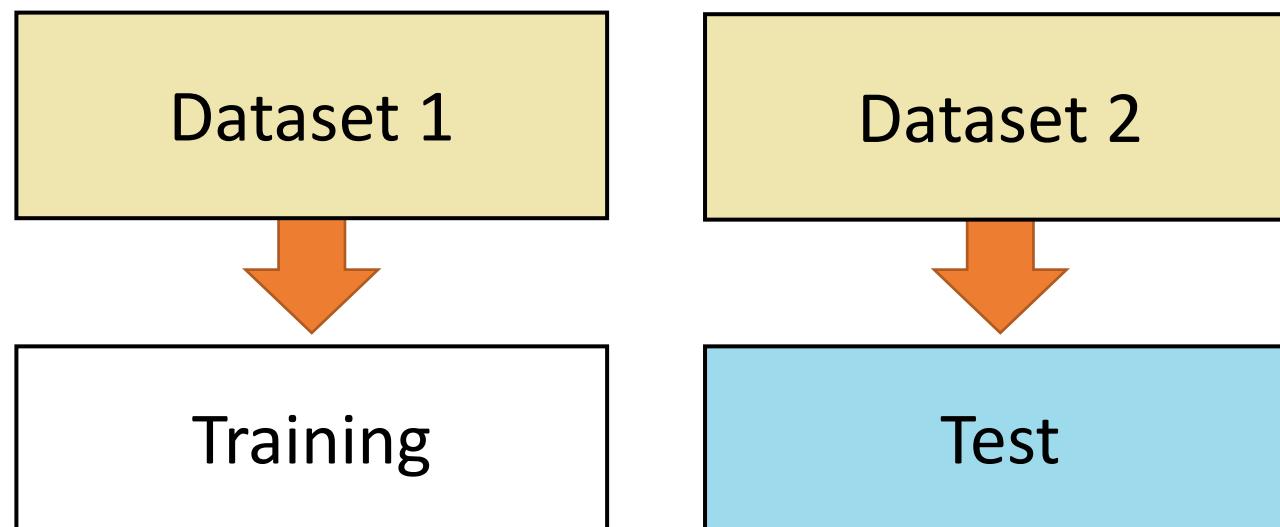
# Performance depends on dataset complexity

$$\text{Complexity} = \text{mean} \left( \max_{\forall i, i \neq j} \text{corr} \left( \text{avg}_{C_i}, \text{avg}_{C_j} \right) \right)$$

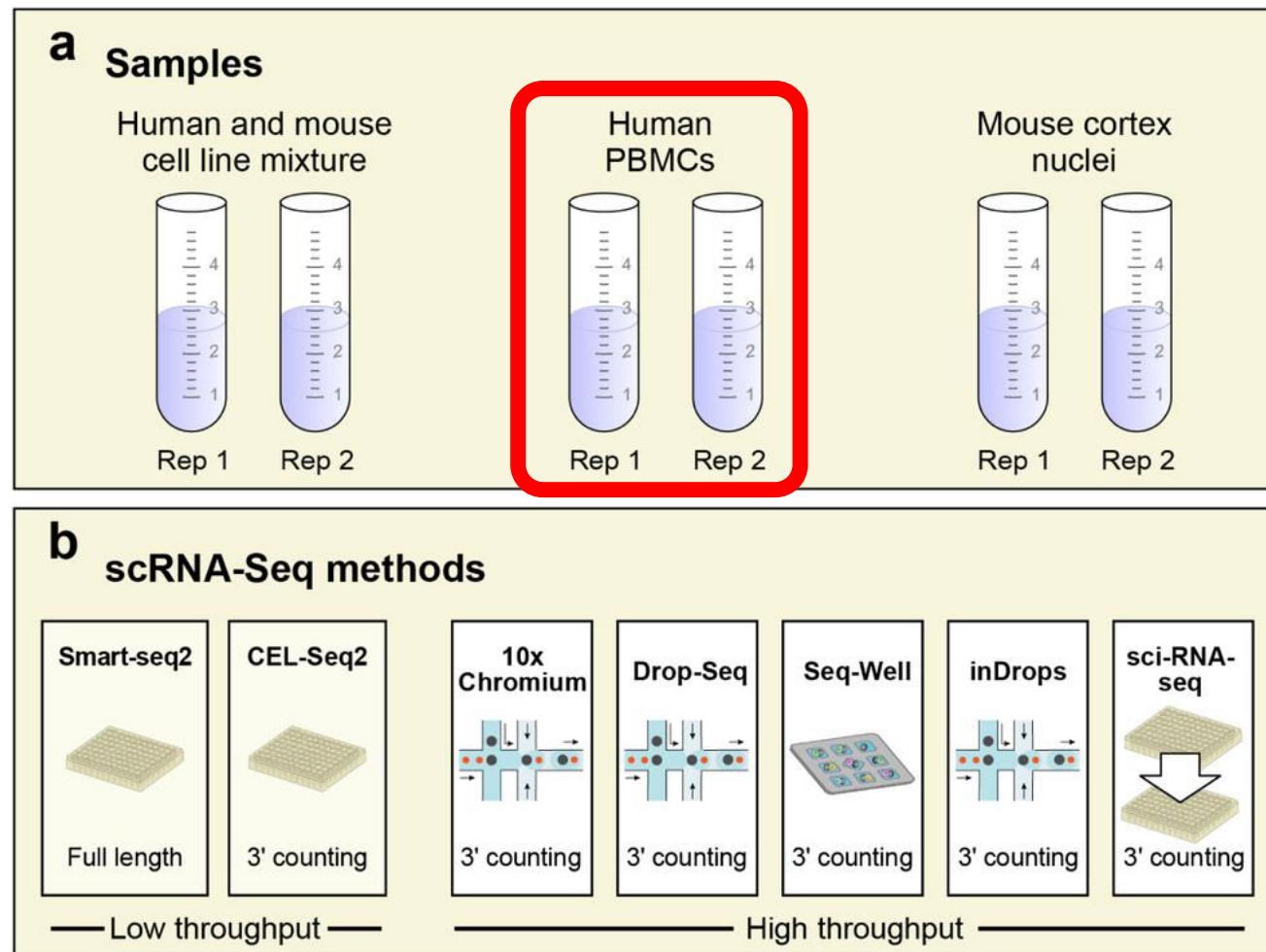


# Experiment 2: inter-dataset evaluation

- Train on one dataset, evaluate on another
- More realistic scenario
- More challenging, data is not aligned



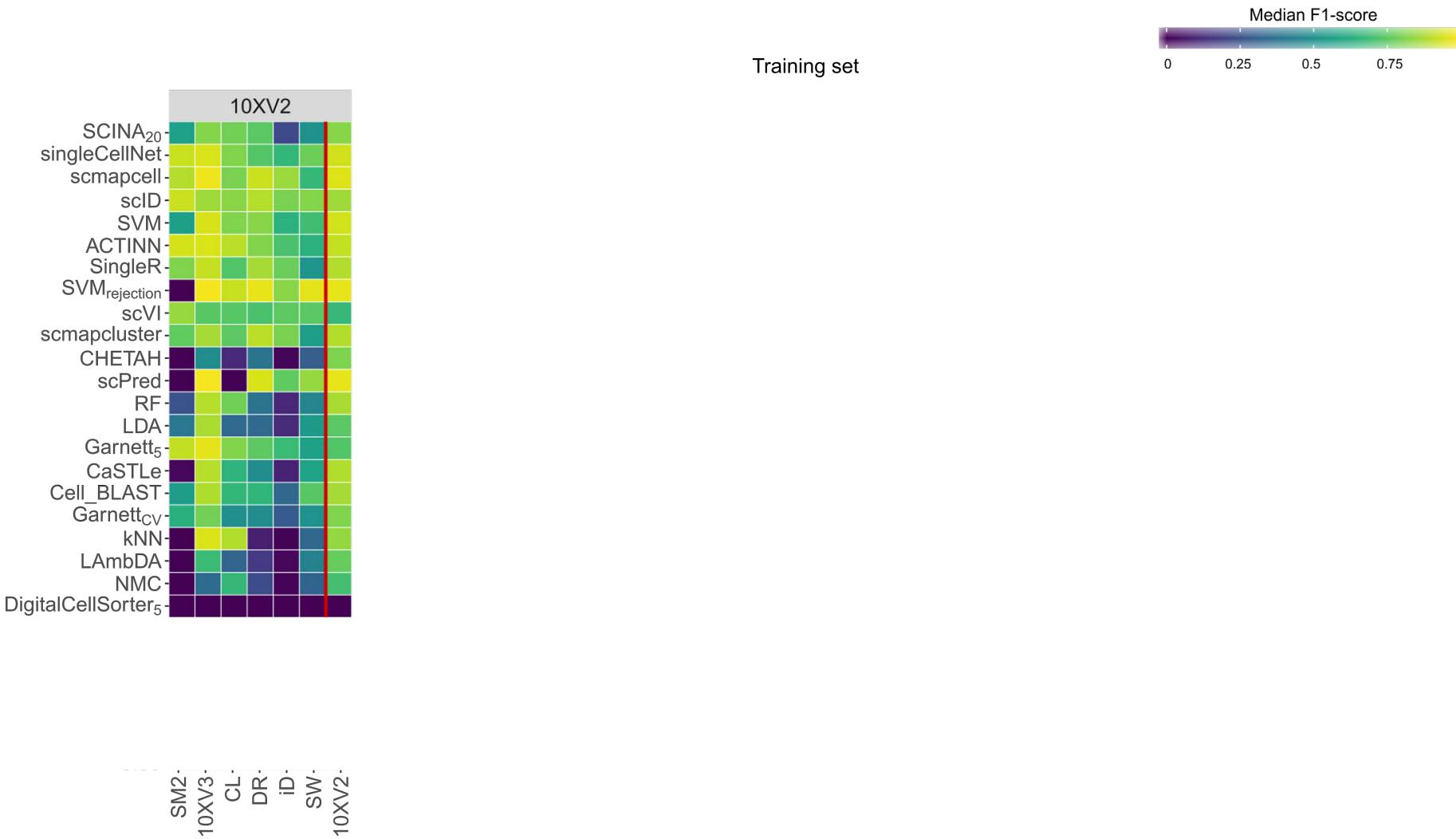
# Experiment 2: inter-dataset evaluation



# Prediction across protocols



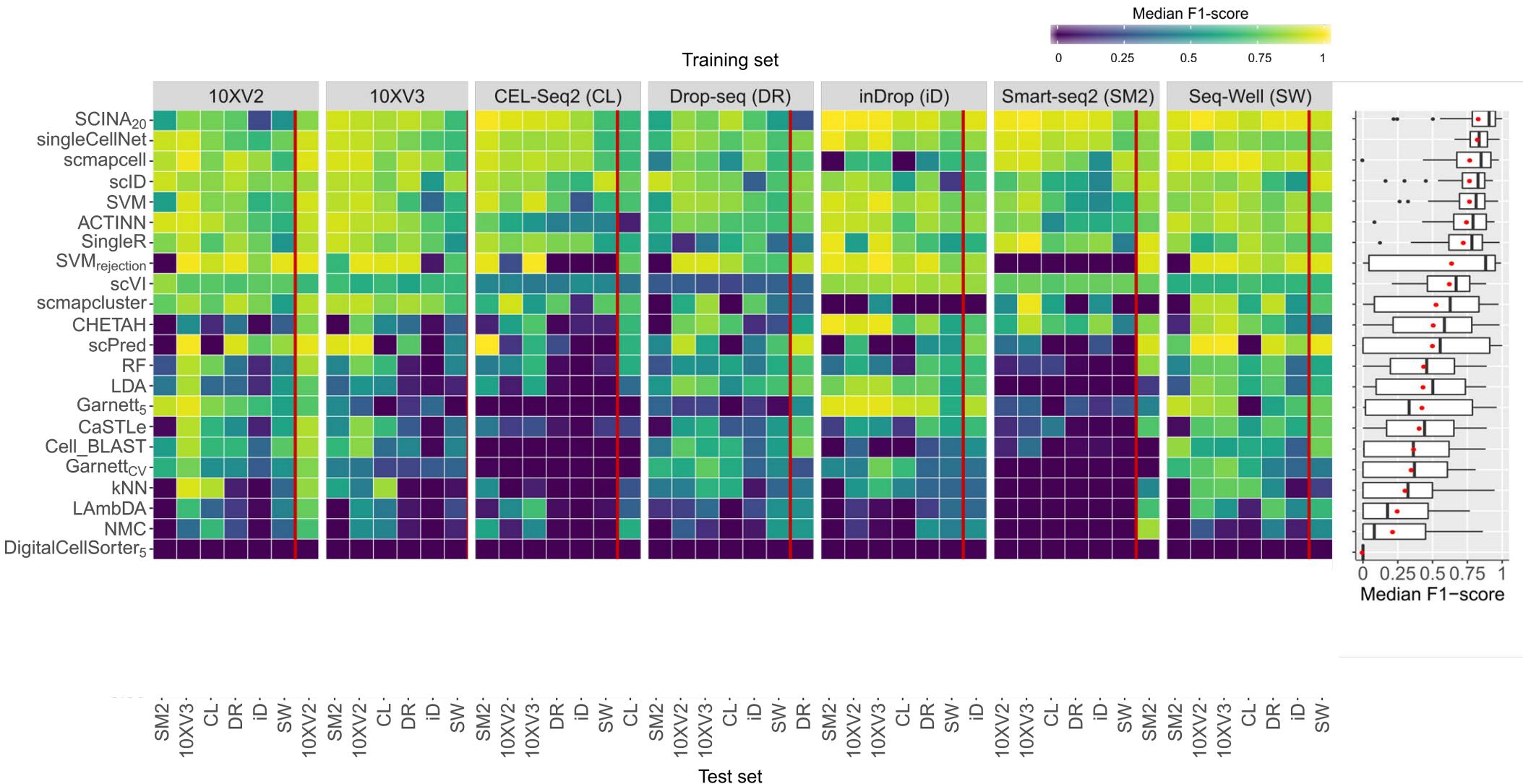
# Prediction across protocols



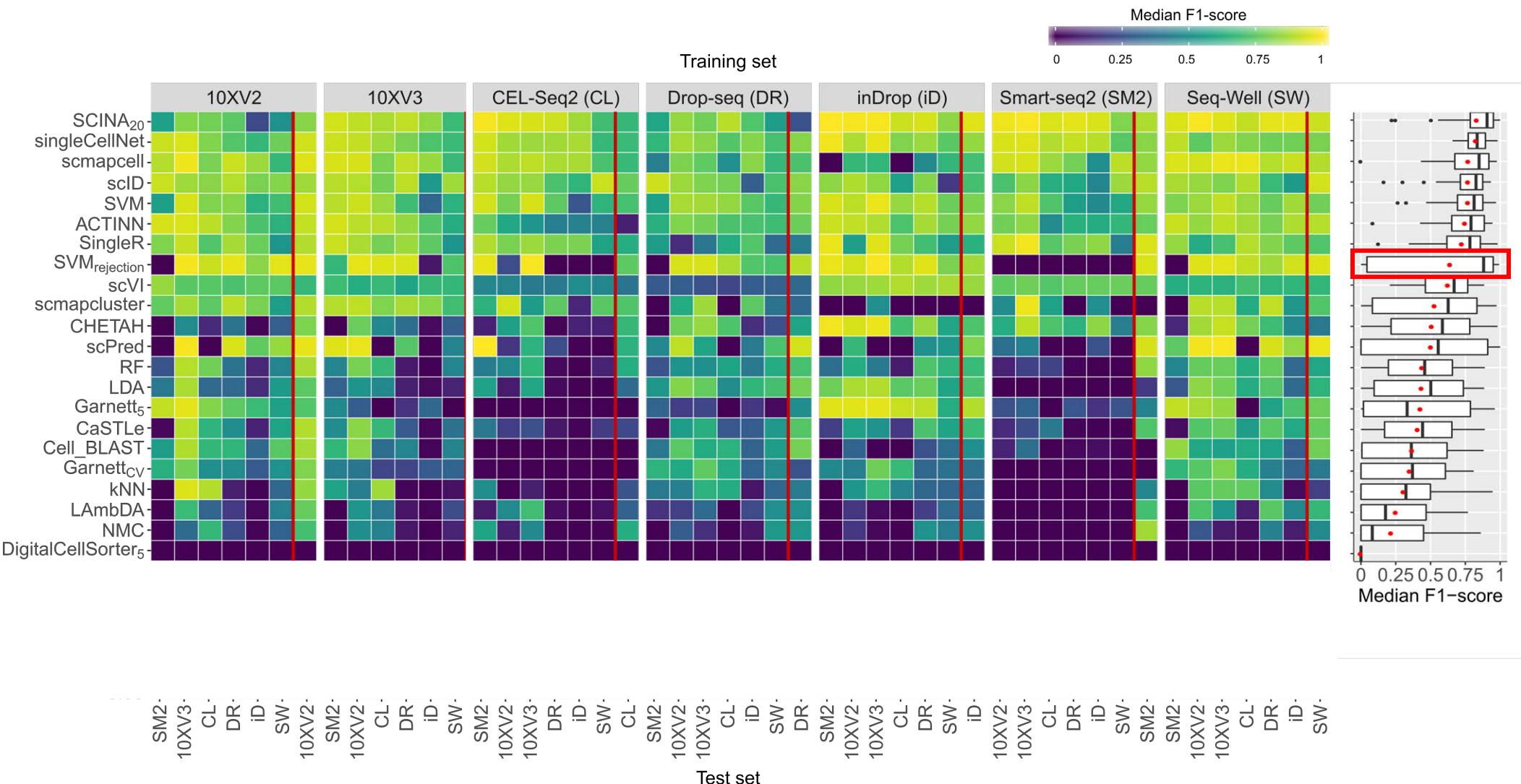
# Prediction across protocols



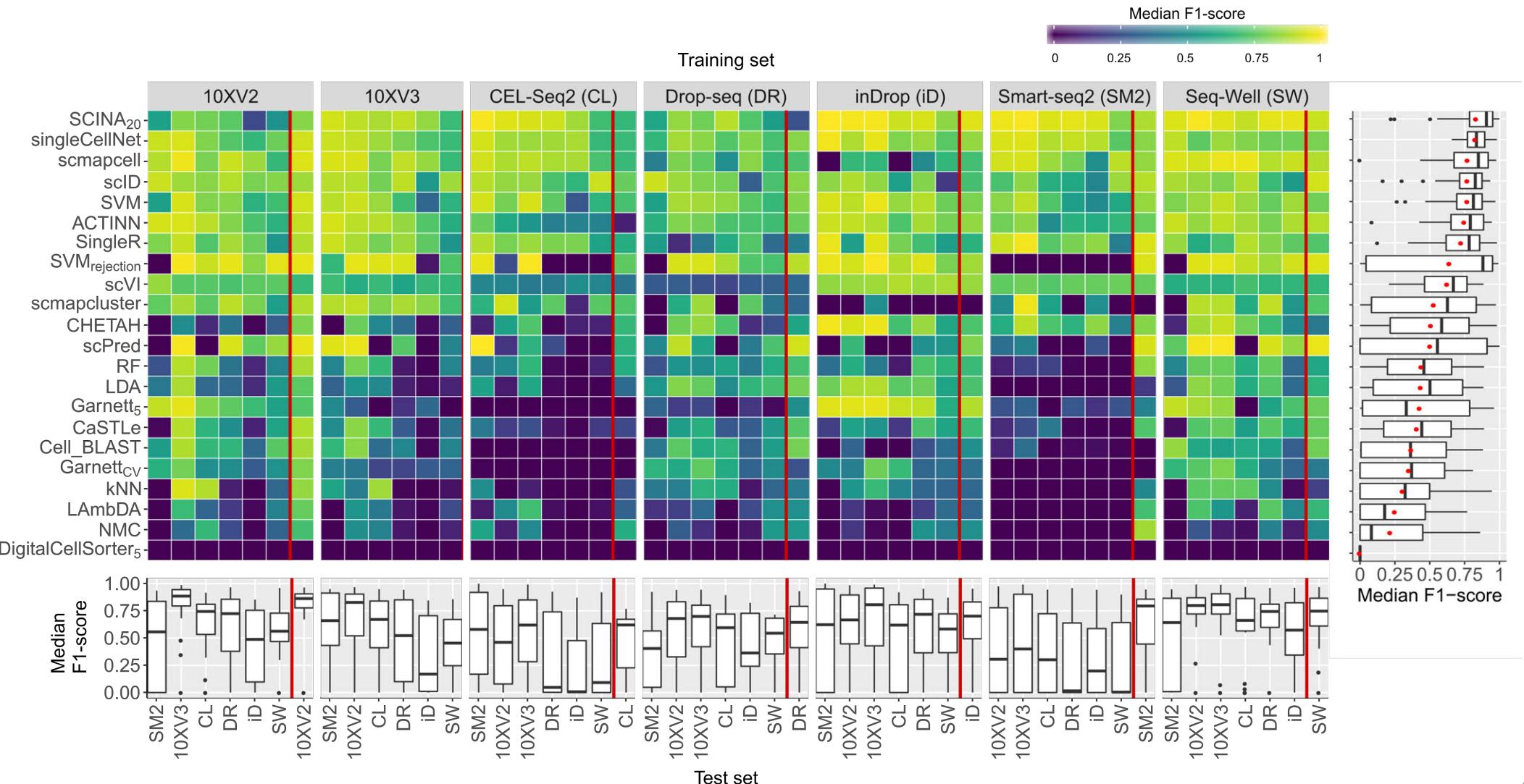
# Prediction across protocols



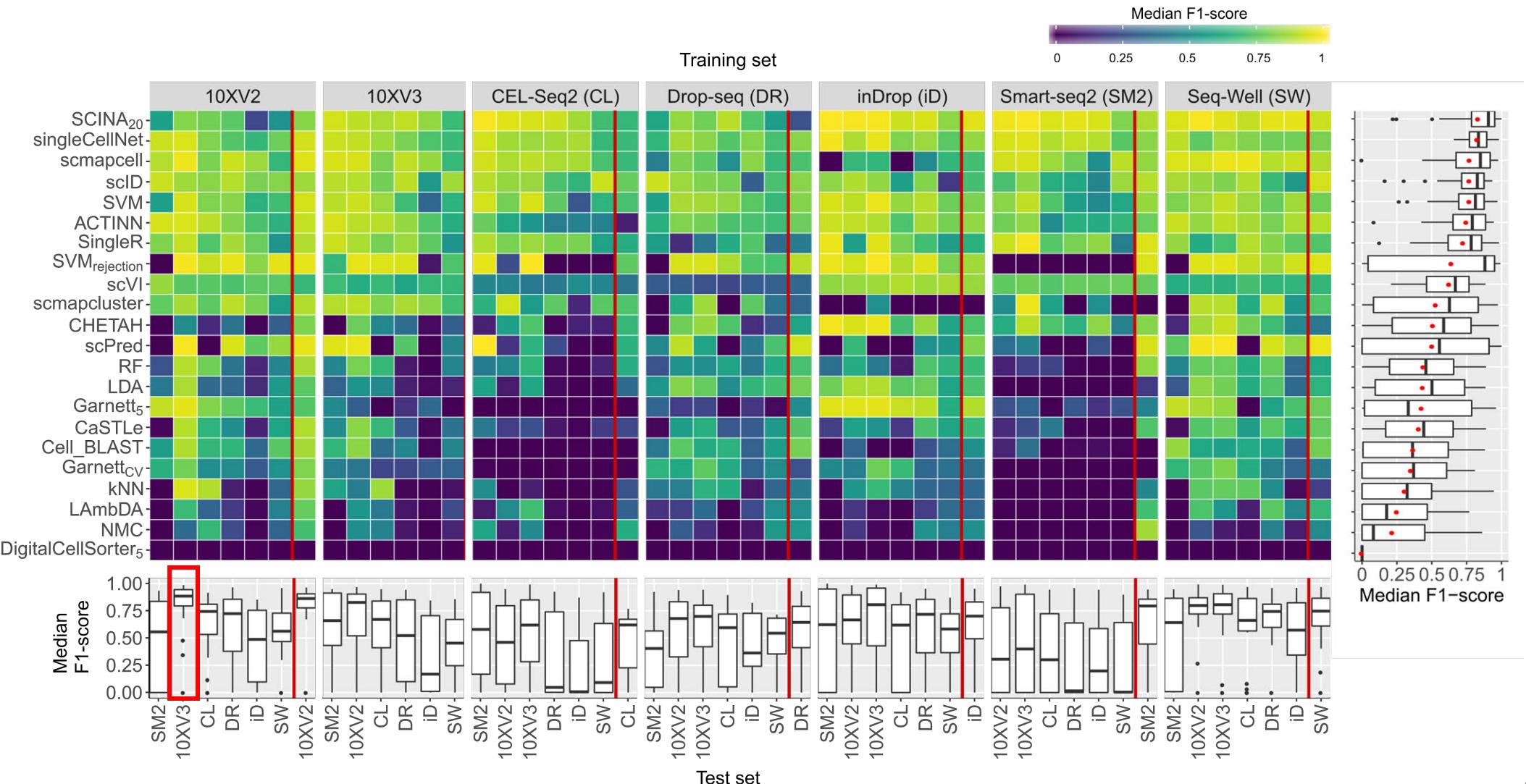
# Prediction across protocols



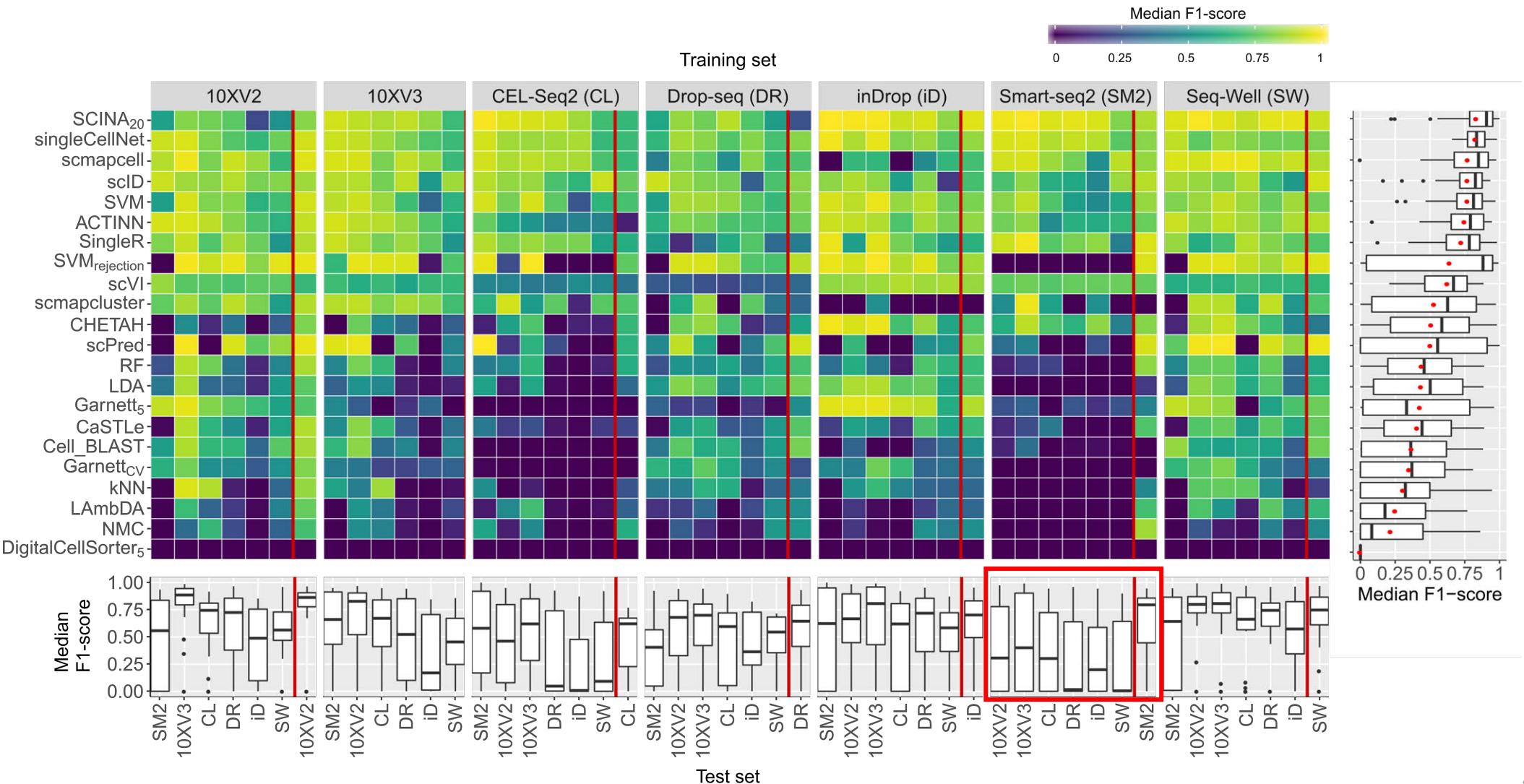
# Prediction across protocols



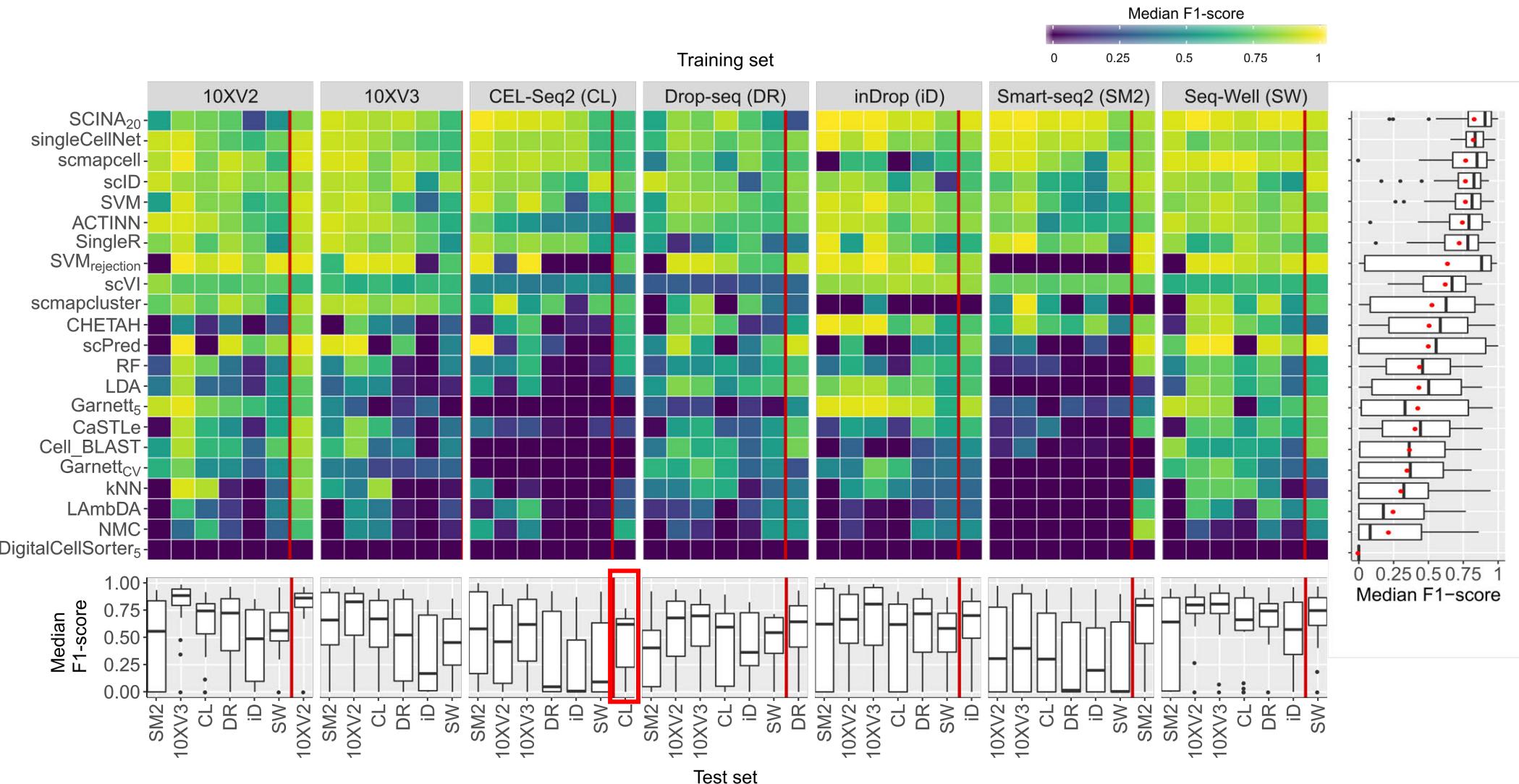
# Prediction across protocols



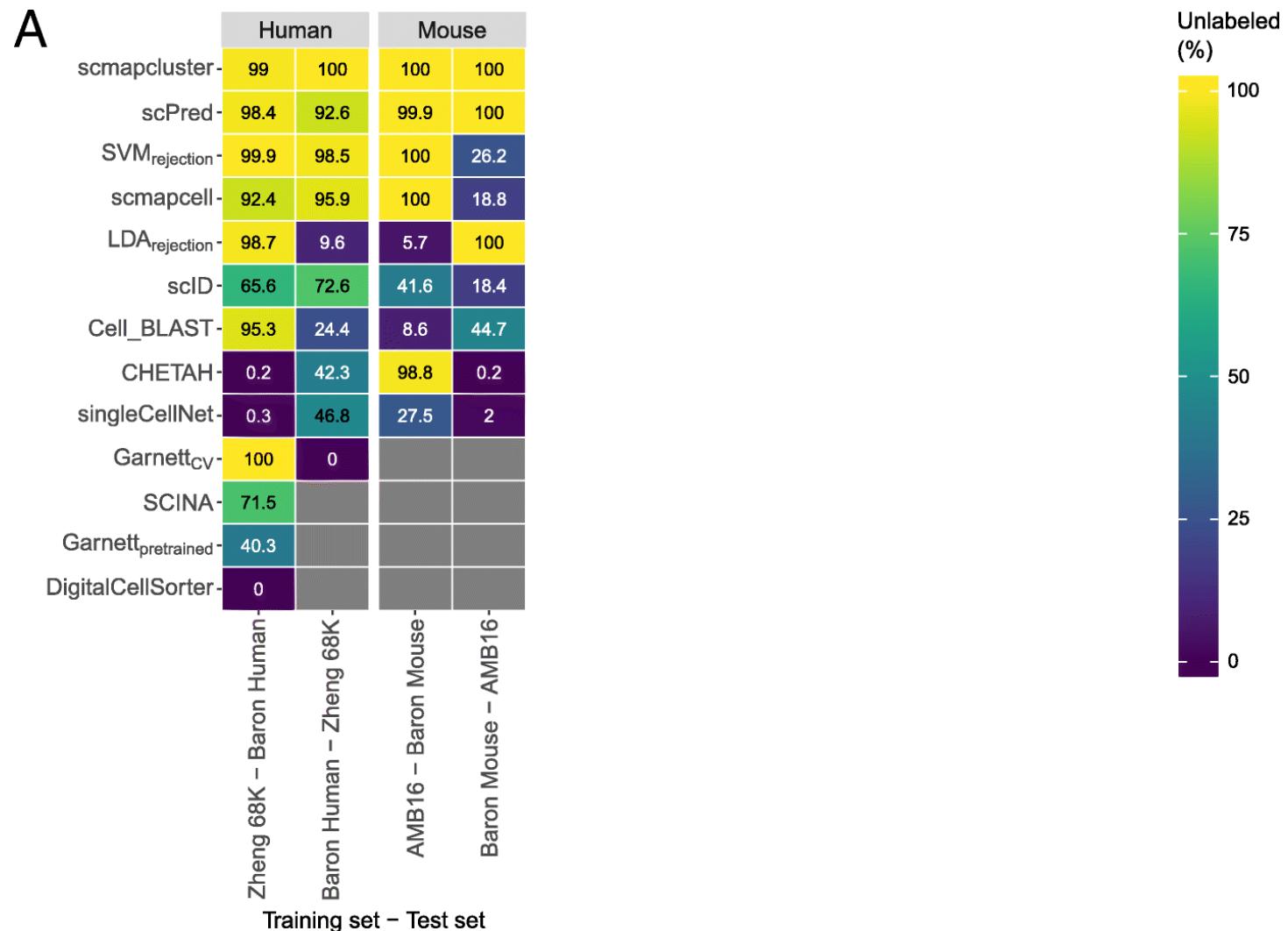
# Prediction across protocols



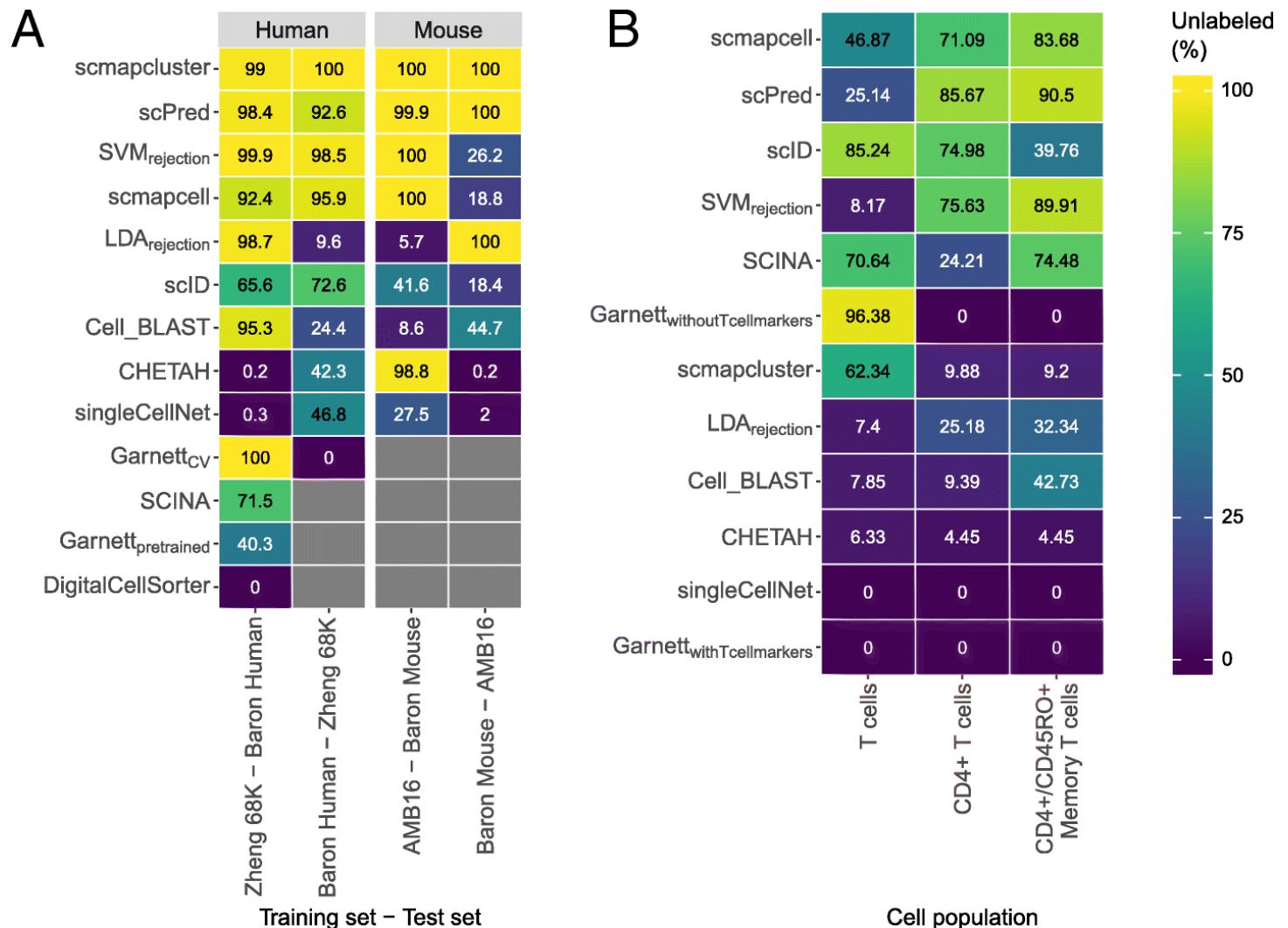
# Prediction across protocols



# Experiment 3: rejection evaluation



# Experiment 3: rejection evaluation



# Performance Summary

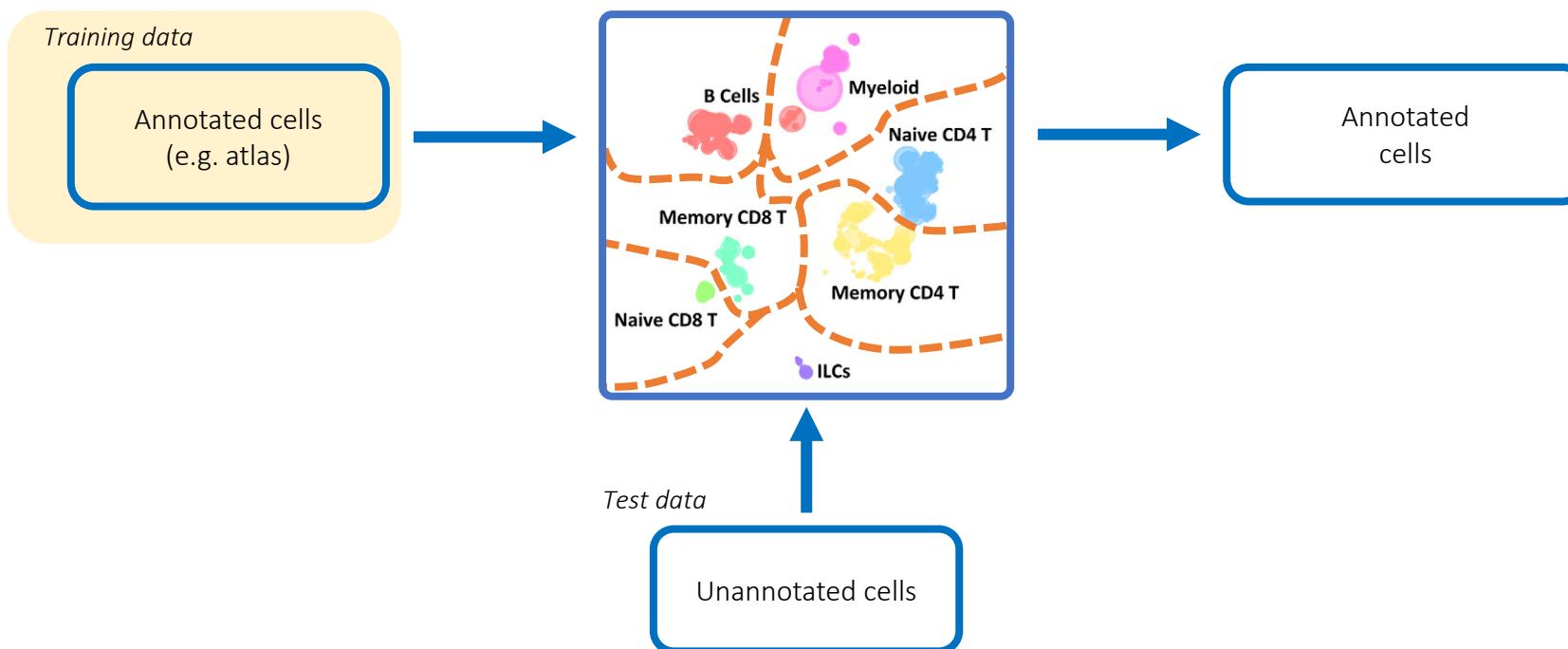


# Conclusions so far

- Simple, off-the-shelf classifiers outperform dedicated single cell methods (see also Köhler et al. bioRxiv 2019)
- Prior-knowledge does not improve performance (highly dependent on selected markers)
- Rejection is difficult
- SnakeMake pipeline:  
[https://github.com/tabdelaal/scRNAseq\\_Benchmark/](https://github.com/tabdelaal/scRNAseq_Benchmark/)

# Still, challenges remain

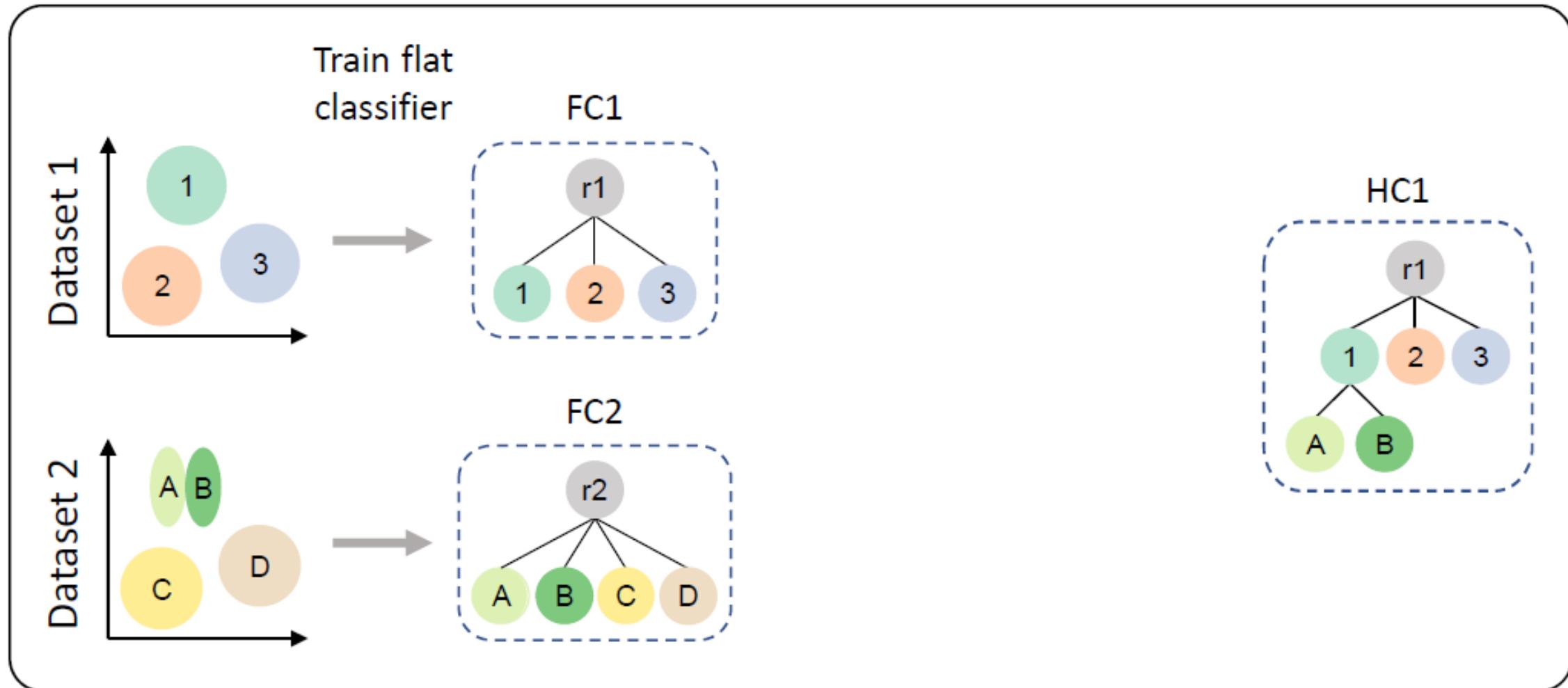
- Incomplete/missing reference atlas
- Inconsistent labels across datasets
- Sharing data is an issue (scArches, Lotfollahi et al.; bioRxiv 2020)



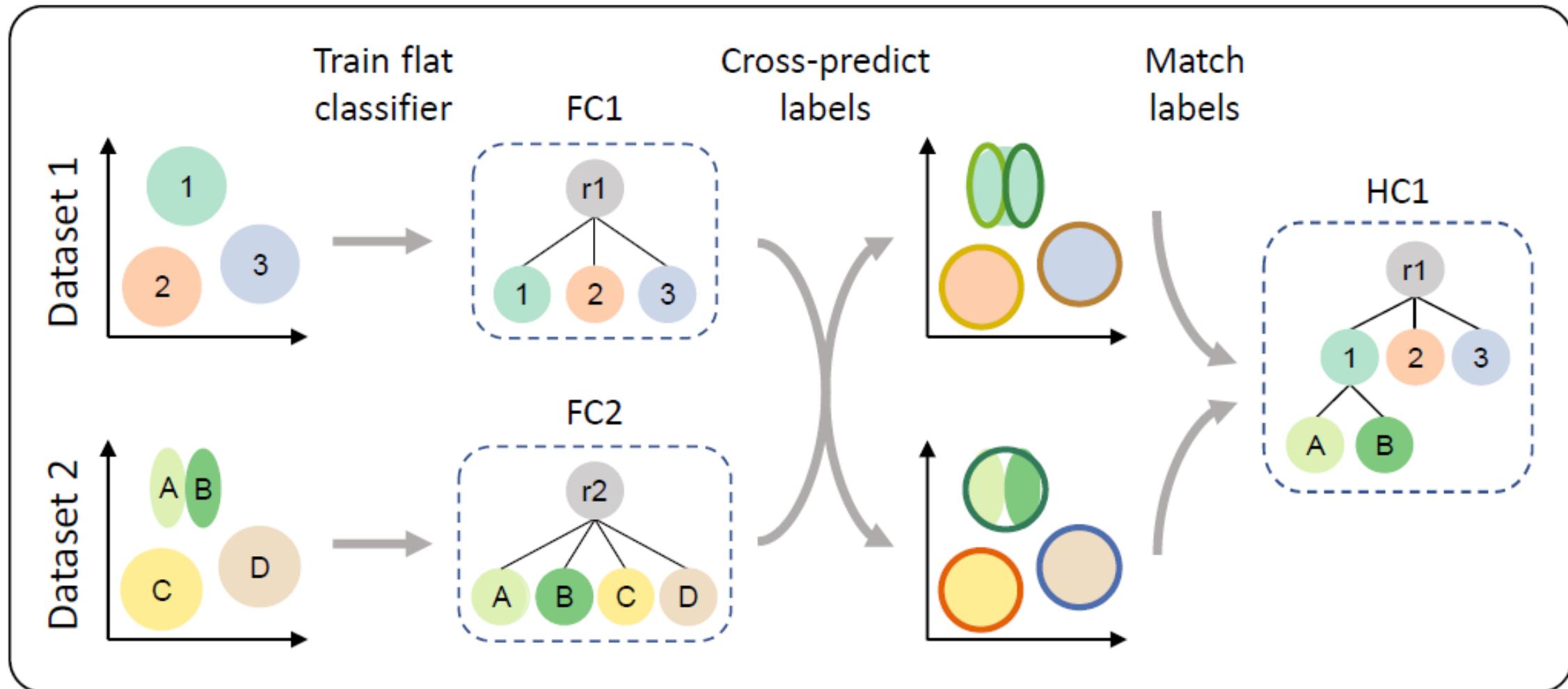
# Hierarchical Progressive Learning



# Hierarchical Progressive Learning



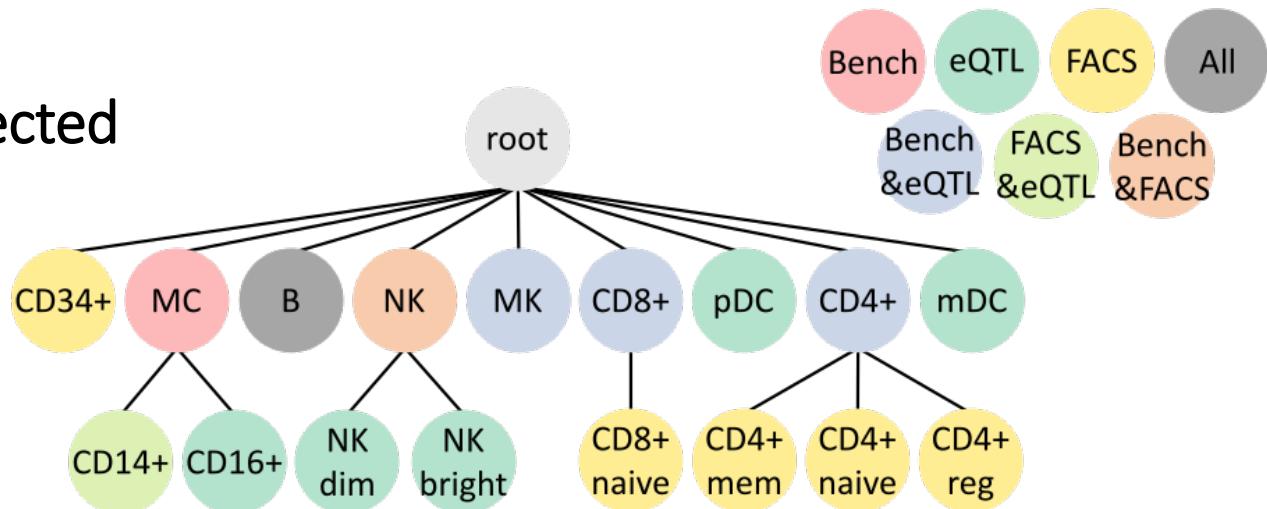
# Hierarchical Progressive Learning



# Tree construction

Cell population	Batch1 eQTL	Batch2 Bench 10Xv2	Batch3 FACS
CD19+ B	812	676	2,000
Monocytes (MC)		1,194	
CD14+	2,081		2,000
CD16+	274		
CD4+ T	13,523	1,458	
Reg.			2,000
Naive			2,000
Memory			2,000
CD8+ T	4,195	2,128	
Naive			2,000
Megakaryocyte (MK)	142	433	
NK cell		429	2,000
CD56+ bright	355		
CD56+ dim	2,415		
Dendritic			
Plasmacytoid (pDC)	101		
Myeloid (mDC)	455		
CD34+		2,000	

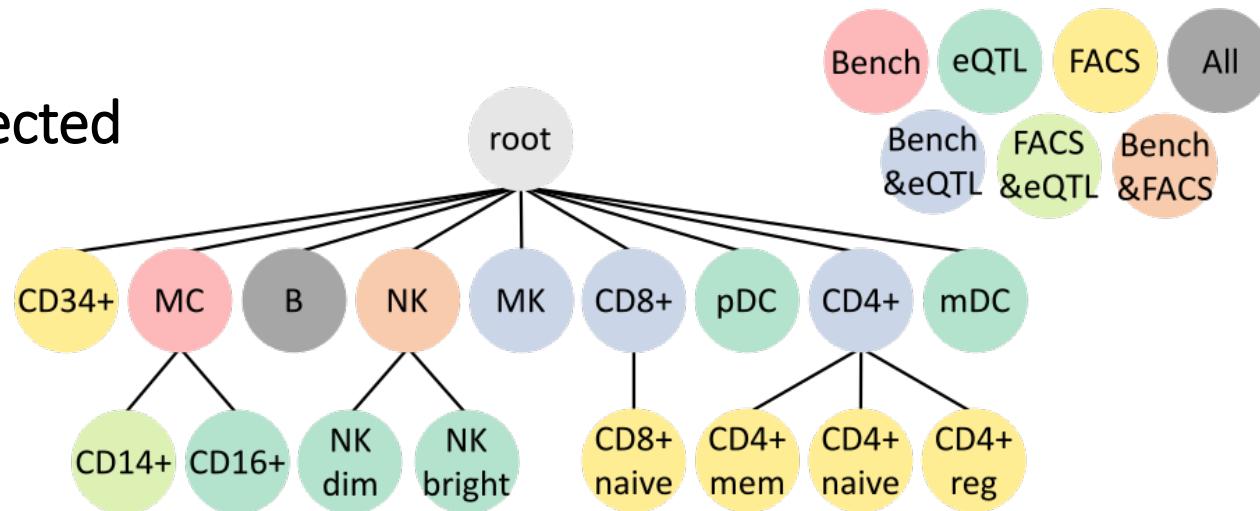
Expected



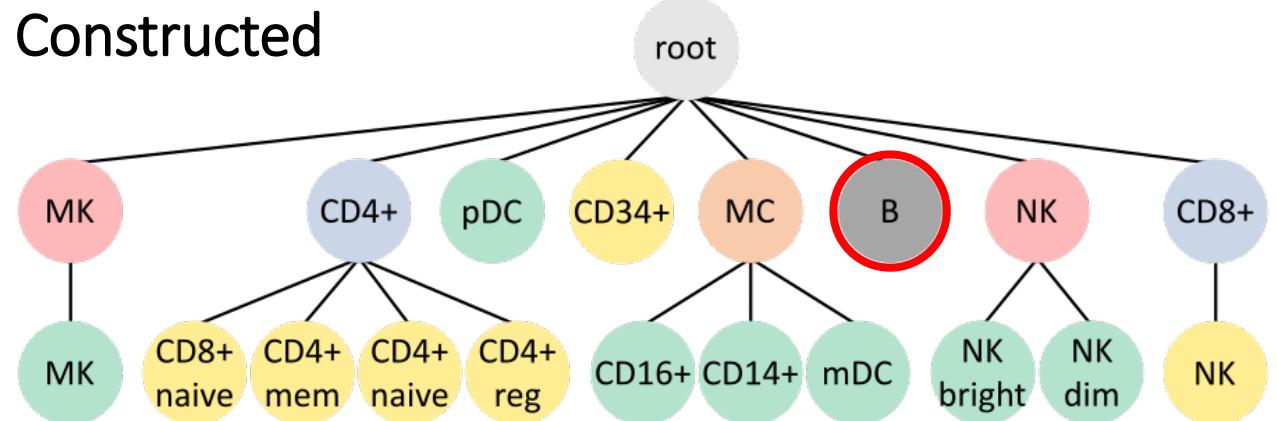
# Tree construction

Cell population	Batch1 eQTL	Batch2 Bench 10Xv2	Batch3 FACS
CD19+ B	812	676	2,000
Monocytes (MC)		1,194	
CD14+	2,081		2,000
CD16+	274		
CD4+ T	13,523	1,458	
Reg.			2,000
Naive			2,000
Memory			2,000
CD8+ T	4,195	2,128	
Naive			2,000
Megakaryocyte (MK)	142	433	
NK cell		429	2,000
CD56+ bright	355		
CD56+ dim	2,415		
Dendritic			
Plasmacytoid (pDC)	101		
Myeloid (mDC)	455		
CD34+			2,000

Expected



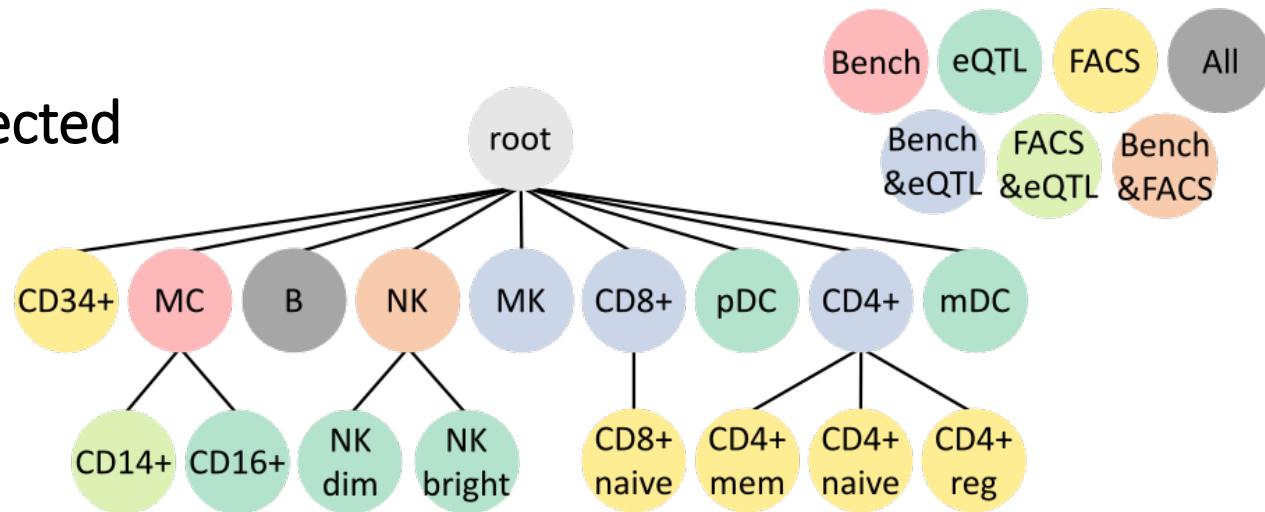
Constructed



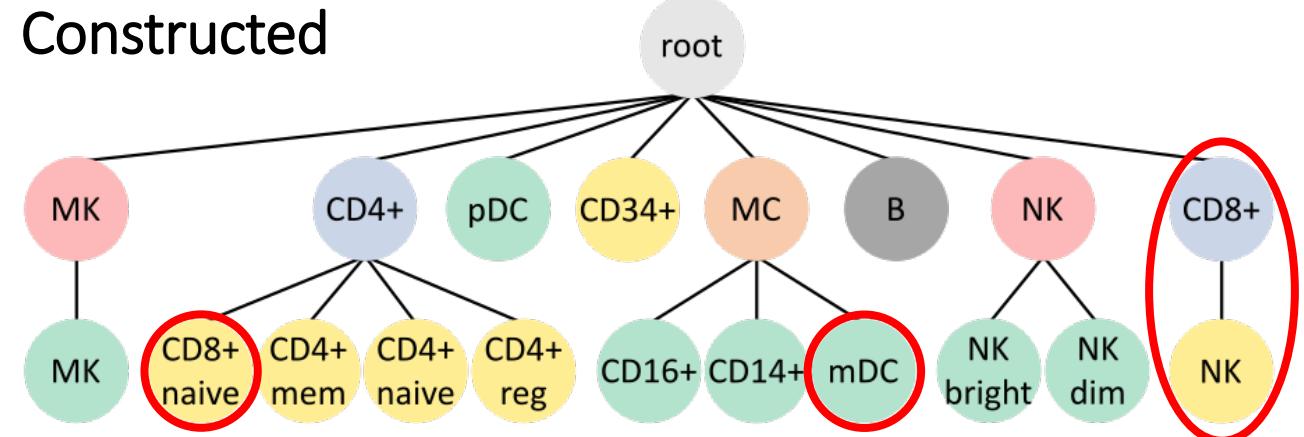
# Tree construction

Cell population	Batch1 eQTL	Batch2 Bench 10Xv2	Batch3 FACS
CD19+ B	812	676	2,000
Monocytes (MC)		1,194	
CD14+	2,081		2,000
CD16+	274		
CD4+ T	13,523	1,458	
Reg.			2,000
Naive			2,000
Memory			2,000
CD8+ T	4,195	2,128	
Naive			2,000
Megakaryocyte (MK)	142	433	
NK cell		429	2,000
CD56+ bright	355		
CD56+ dim	2,415		
Dendritic			
Plasmacytoid (pDC)	101		
Myeloid (mDC)	455		
CD34+			2,000

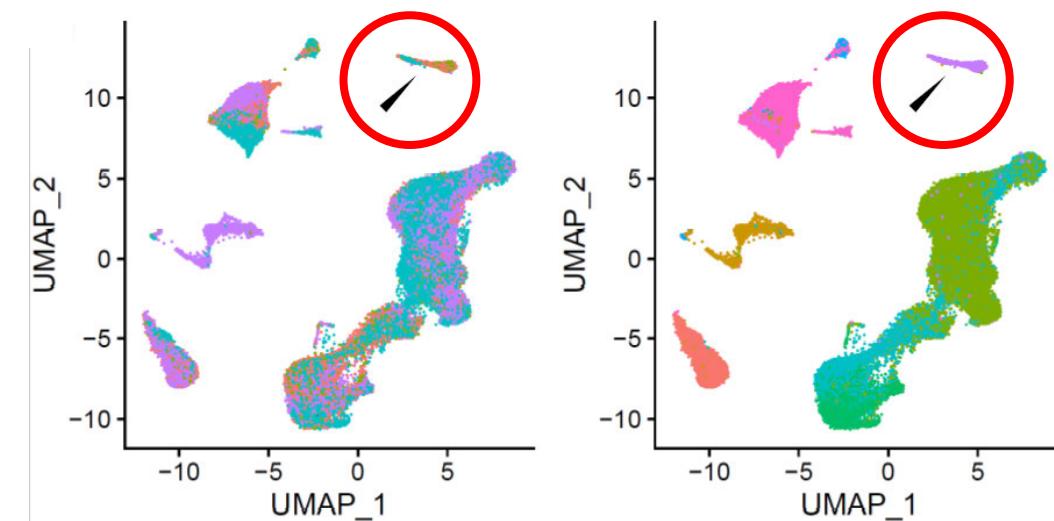
Expected



Constructed



# Tree construction



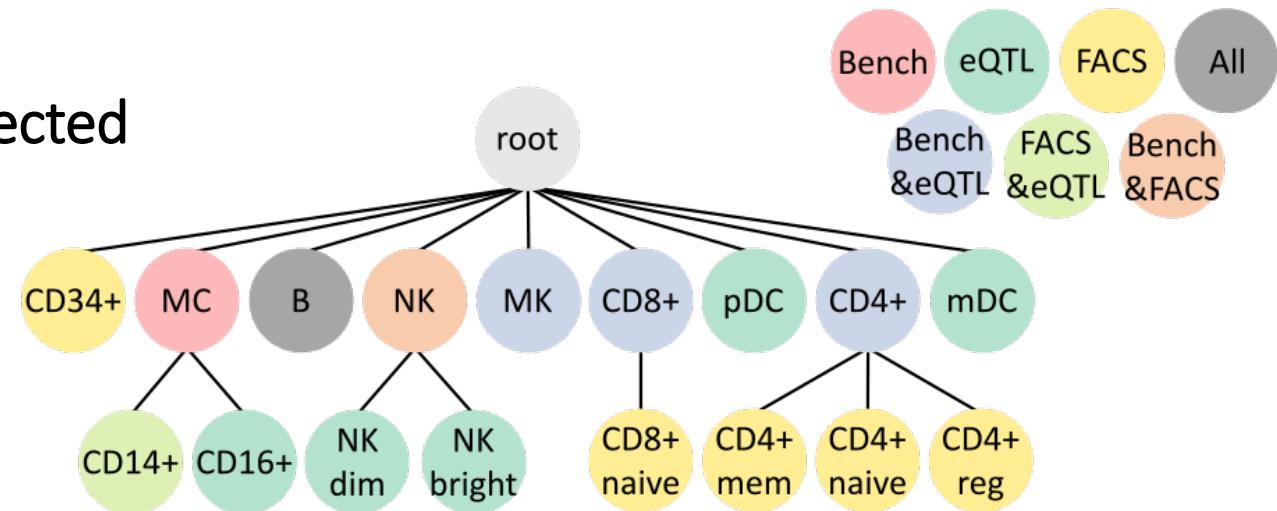
## Dataset

- 10Xv2
- 10Xv3
- eQTL
- FACS

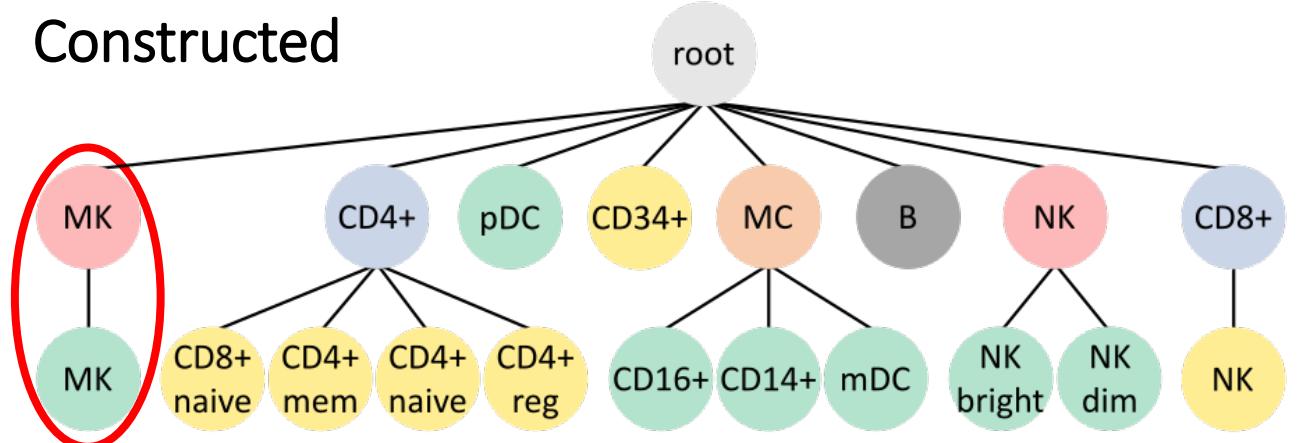
## Cell population

- B cell
- CD34+
- CD4+ T
- CD56+ NK
- CD8+ T
- Dendritic cell
- Megakaryocyte
- Monocytes

## Expected

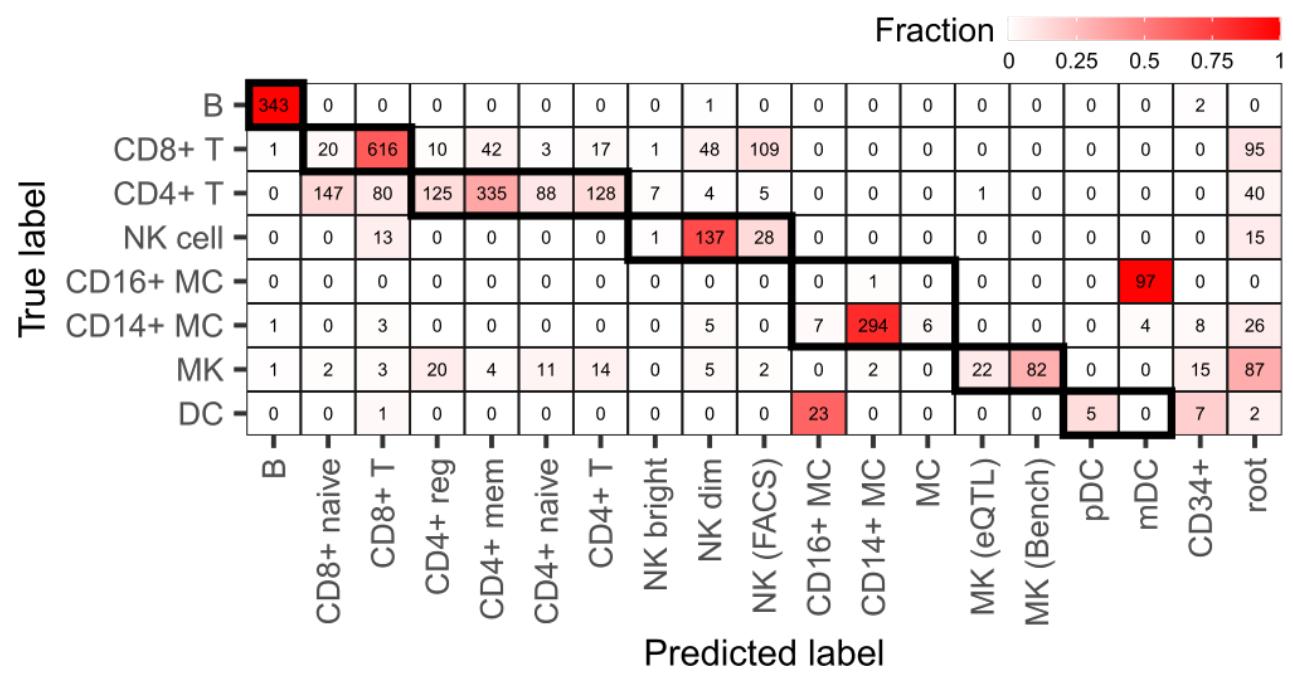


## Constructed



# Classification performance

Cell population	Batch1 eQTL	Batch2 Bench 10Xv2	Batch3 FACS	Testset Bench 10Xv3
CD19+ B	812	676	2,000	346
Monocytes (MC)		1,194		
CD14+	2,081		2,000	354
CD16+	274			98
CD4+ T	13,523	1,458		960
Reg.			2,000	
Naive			2,000	
Memory			2,000	
CD8+ T	4,195	2,128		962
Naive			2,000	
Megakaryocyte (MK)	142	433		270
NK cell		429	2,000	194
CD56+ bright	355			
CD56+ dim	2,415			
Dendritic				38
Plasmacytoid (pDC)	101			
Myeloid (mDC)	455			
CD34+		2,000		



# Summary

- Cell identification is moving from unsupervised (clustering/visualization) to supervised (classification) learning
- Comprehensive benchmark of classifiers for single-cell RNA-seq data helps both users and developers
- Continuous learning from a growing reference atlas by combining multiple annotated datasets into a hierarchical classifier (scHPL)

# Thank You!



 a.mahfouz@lumc.nl  
 mahfouzlab.org  
 @ahmedElkoussy



CHAN  
ZUCKERBERG  
INITIATIVE



IMID



Single-cell eQTL  
Consortium