# Celltype prediction

Åsa Björklund

asa.bjorklund@scilifelab.se

Ahmed Mahfouz
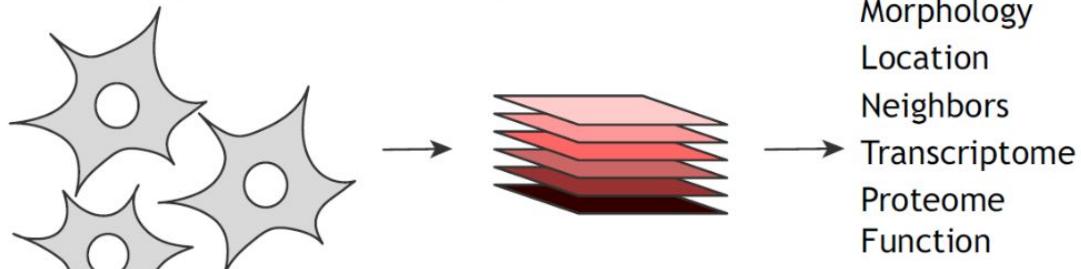Leiden University Medical Center / TU Delft

NBIS
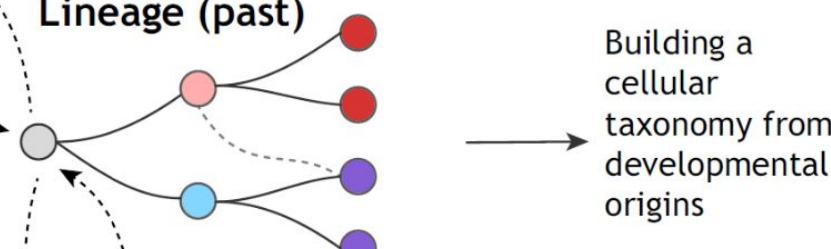NATIONAL BIOINFORMATICS
INFRASTRUCTURE SWEDEN

SciLifeLab

# Outline

- Introduction

- Normalization

- Removal of confounders

- Gene set selection

# Cell identity



**Phenotype and function (present)**

Morphology
Location
Neighbors
Transcriptome
Proteome
Function

**Lineage (past)**

Building a cellular taxonomy from developmental origins

**State (future)**

Distinguishing between cell type and cell state
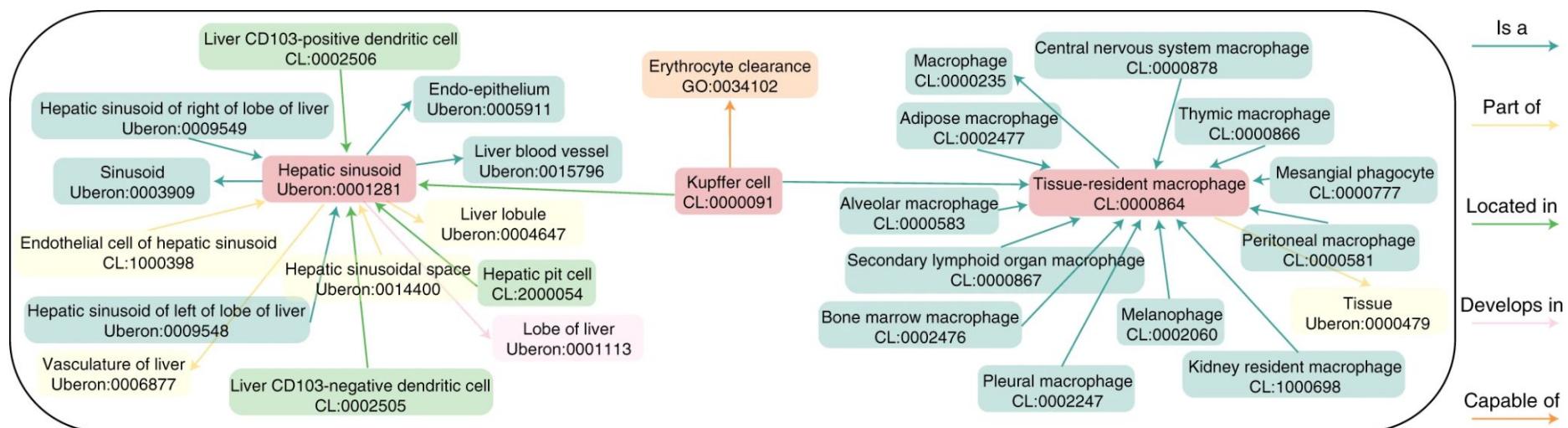
(Morris, *Development* 2019)

# Why do we want to classify celltypes?

- In a novel tissue - what celltypes are there?
- Compare same celltype across conditions.
- Compare abundance of celltypes across conditions.
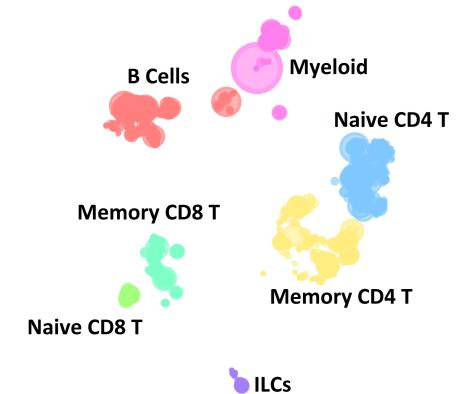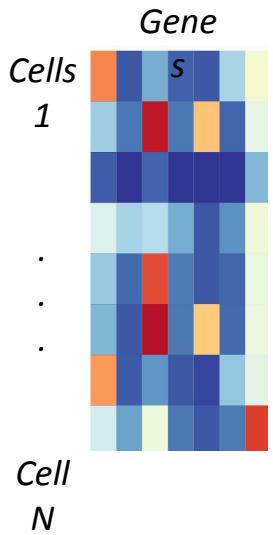- Infer communication between celltypes
- …..

# Celltype ontologies

We need a standardized way of classifying celltypes.
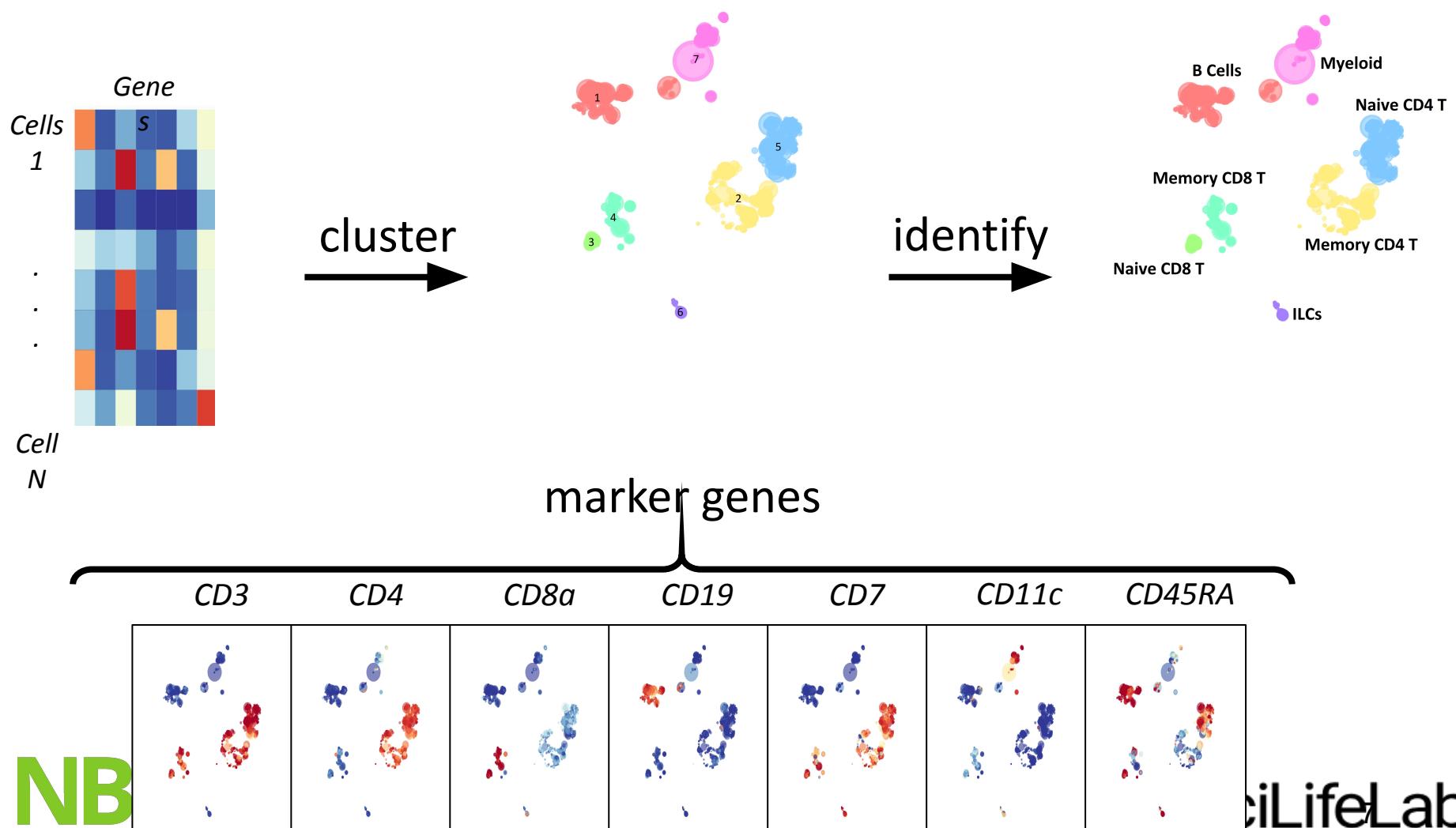Mainly driven by cell atlas projects.

Including HuBMAP, Human Cell Atlas (HCA), cellxgene, Single Cell Expression Atlas, BRAIN Initiative Cell Census Network (BICCN), ArrayExpress, The Cell Image Library, ENCODE, and FANTOM5,



(Osumi-Sutherland, Nature Cell Biol 2021)

# How can we identify cell populations?



*Gene*
*s*

*Cells*
*1*

.
.
.

*Cell*
*N*

B Cells

Myeloid

Naive CD4 T

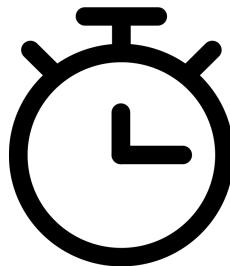Memory CD8 T

Memory CD4 T

Naive CD8 T
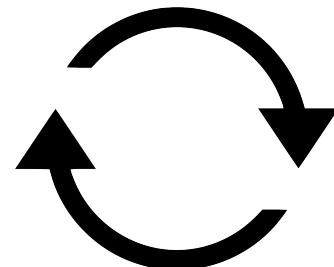
ILCs

# How can we identify cell populations?

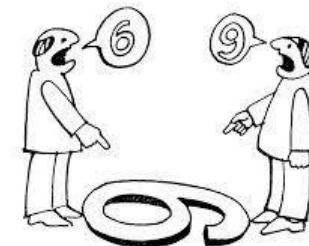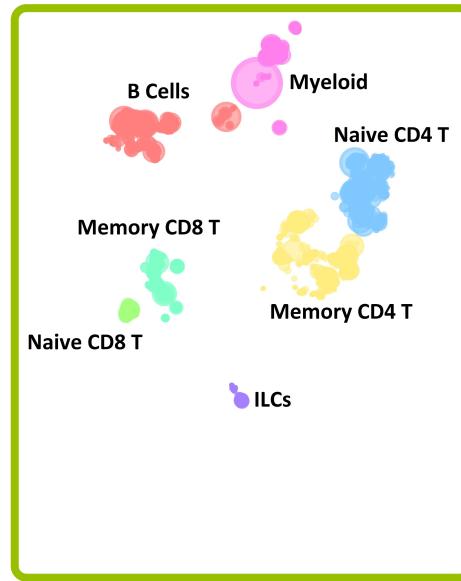# Unsupervised celltype identification is problematic
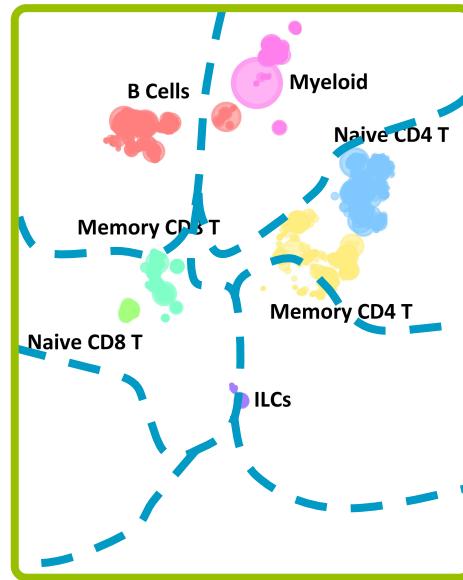
**Time consuming**

**Not reproducible**

**Subjective**
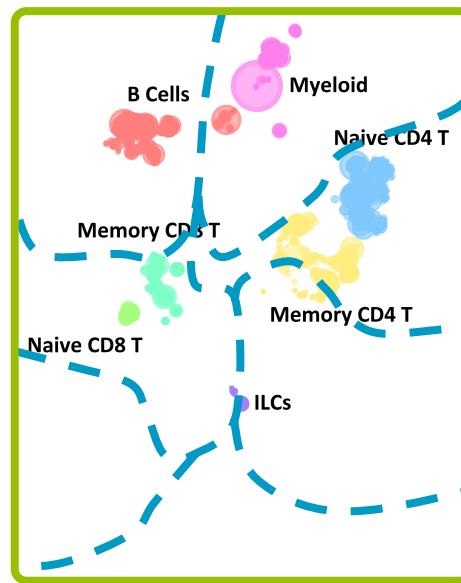
# Can we automatically identify cell populations?

# Can we automatically identify cell populations?

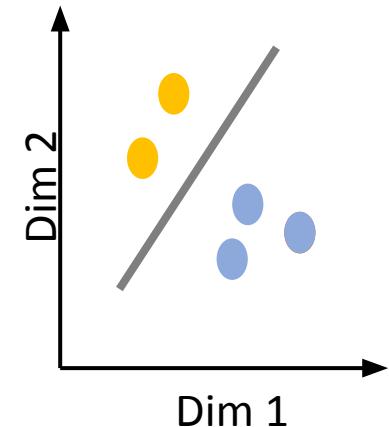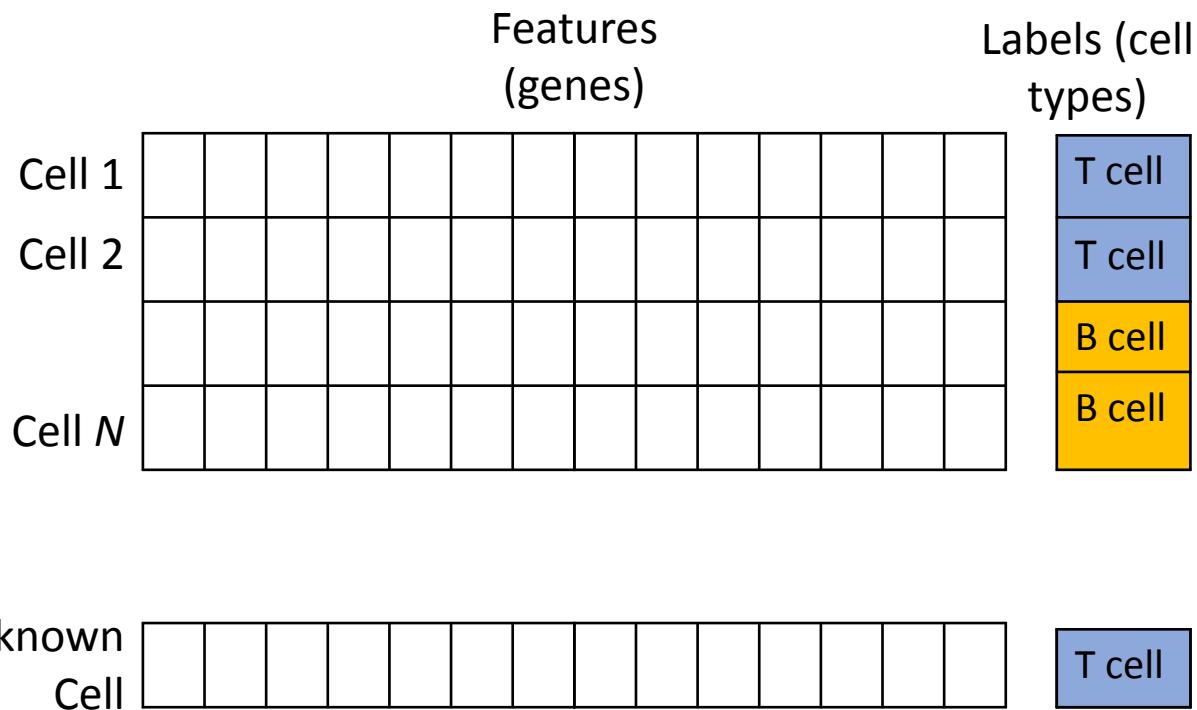# Can we automatically identify cell populations?

## Clustering

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



Myeloid
B Cells
Naive CD4 T
Memory CD8 T
Naive CD8 T
Memory CD4 T
ILCs

## Classification

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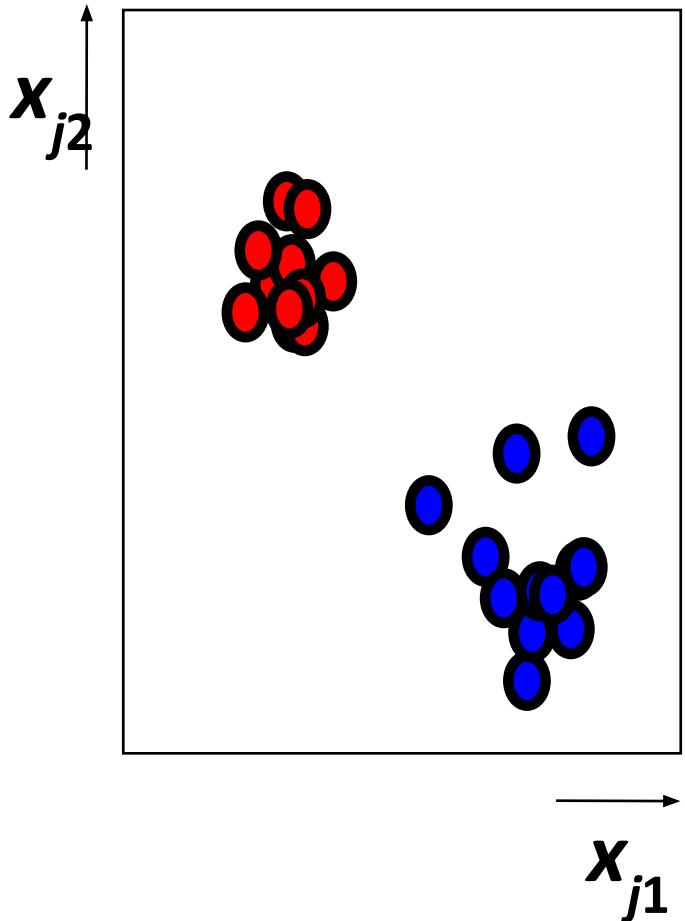
# Classification

Features
(genes)

Labels (cell
types)

| | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
Cell 1

Cell 2

Cell *N*

| T cell |
|---|
| T cell |
| B cell |
| B cell |



Dim 2

Dim 1

Unknown
Cell

| | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

| T cell |
|---|

# Classifier training

- Dataset: for $j^{\text{th}}$ cell:
  - gene expressions $\boldsymbol{x}_j$
  - class label: $y_j \in \{1=T, -1=B\}$

- Classifier:

$$\hat{y}_j = W(x_j)$$

- Errors:

$$E = \text{sum}(E_j) \quad E_j = \begin{cases} 1 & \text{if} \quad \hat{y}_j \neq y_j \\ 0 & \text{if} \quad \hat{y}_j = y_j \end{cases}$$
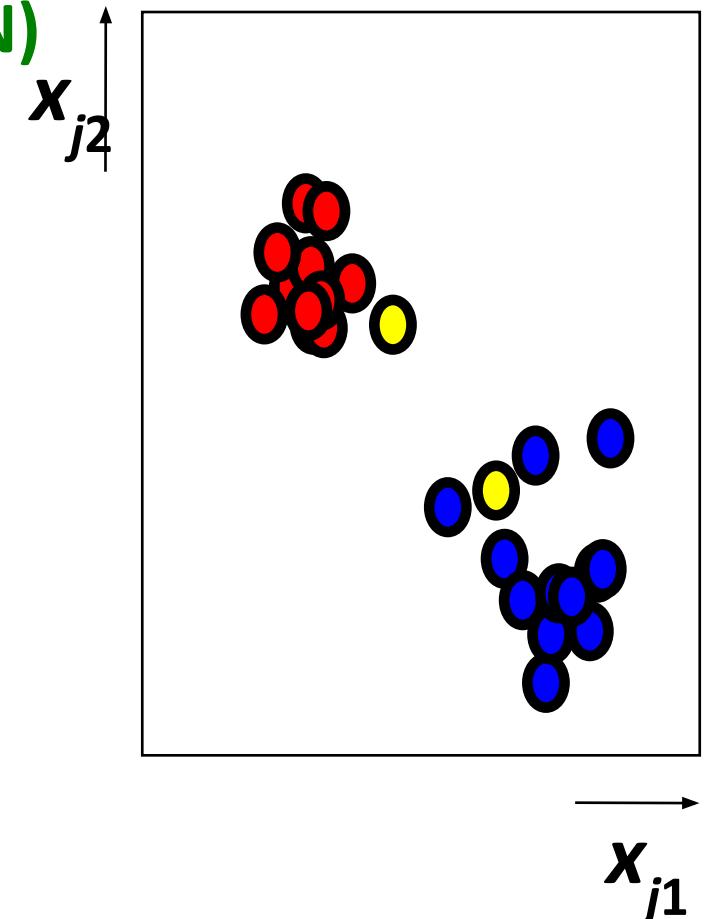
- Place decision boundary (i.e. change $W$) s.t. $E$ is minimal



$\boldsymbol{x}_{j2}$

$\boldsymbol{x}_{j1}$

# Instance Based Learning (Lazy Classification)

- **Example: Nearest neighbor (k-NN)**

$x_{j2}$

- Keep the whole training dataset

- A query example (vector) comes

- Find closest example(s)
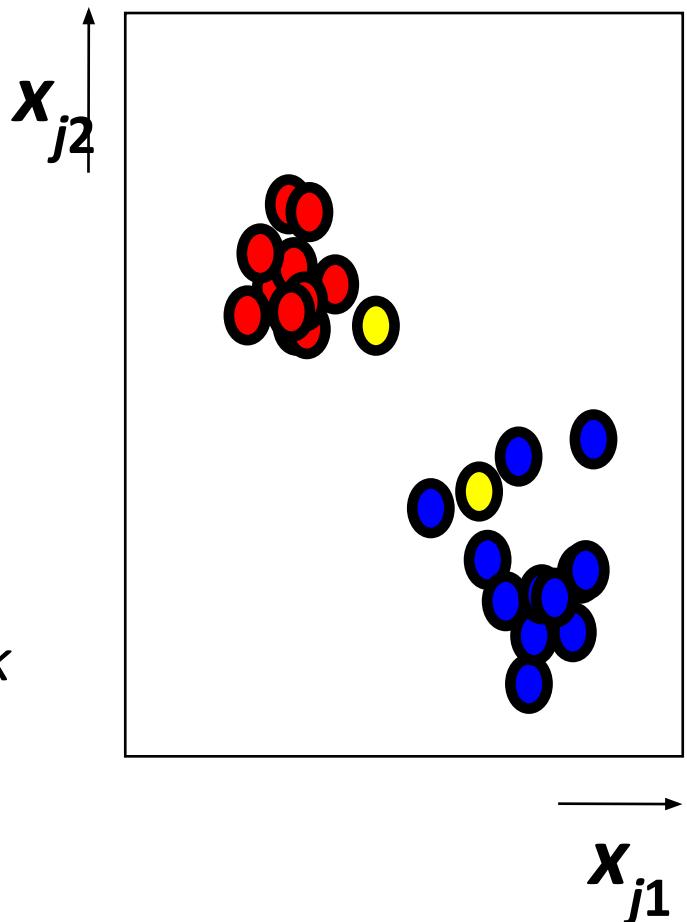
- Predict

- *No actual training*

$x_{j1}$

# Nearest Neighbor (k-NN)

- To make Nearest Neighbor work we need 4 things:

1) Distance metric:

2) How many neighbors to look at?

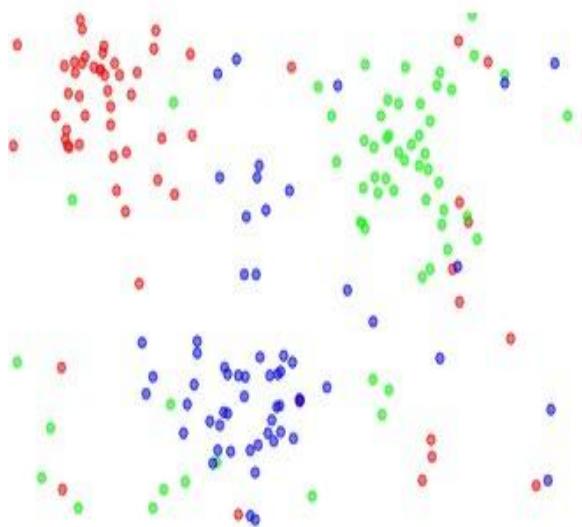3) Weighting function (optional)

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?
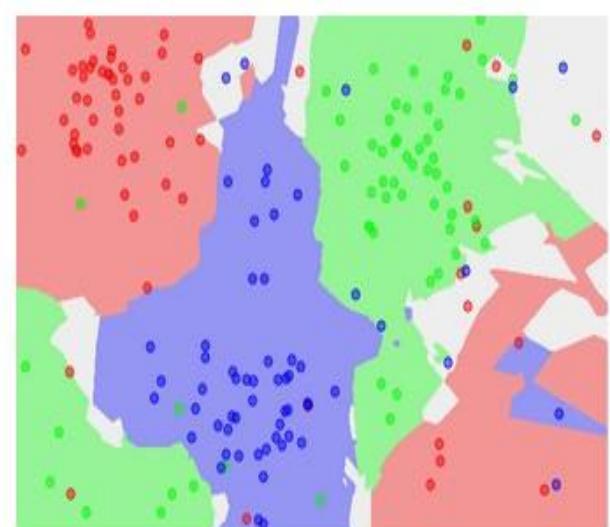  - Predict the average output among $k$ nearest neighbors



$x_{j2}$

$x_{j1}$

# 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 a query q are weighted more strongly. $K_W$: kernel width

- How to fit with the local points?
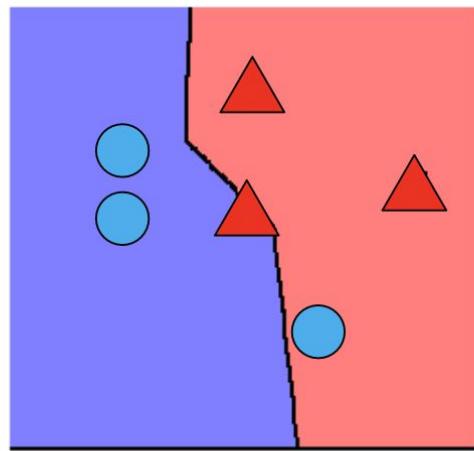  - Predict the weighted average $\dfrac{\sum_i w_i y_i}{\sum_i w_i}$

$w_i$

$d(x_i, q) = 0$

Comparison: K=1, K=2, kernel

# Seurat data transfer

```
pancreas.anchors <- FindTransferAnchors(reference = pancreas.ref, query = pancreas.query, dims =
1:30,
    reference.reduction = "pca")
predictions <- TransferData(anchorset = pancreas.anchors, refdata = pancreas.ref$celltype, dims =
1:30)
pancreas.query <- AddMetaData(pancreas.query, metadata = predictions)
```

```
TransferData(
  anchorset,
  refdata,
  reference = NULL,
  query = NULL,
  query.assay = NULL,
  weight.reduction = "pcaproject",
  l2.norm = FALSE,
  dims = NULL,
  k.weight = 50,
  sd.weight = 1,
  eps = 0,
  n.trees = 50,
  verbose = TRUE,
  slot = "data",
  prediction.assay = FALSE,
  only.weights = FALSE,
  store.weights = TRUE
)
```

# Scanpy data transfer

## scanpy.tl.ingest

```
scanpy.tl.ingest(adata, adata_ref, *, obs=None, embedding_method=
('umap', 'pca'), labeling_method='knn', neighbors_key=None,
inplace=True, **kwargs)                                          [source]
```
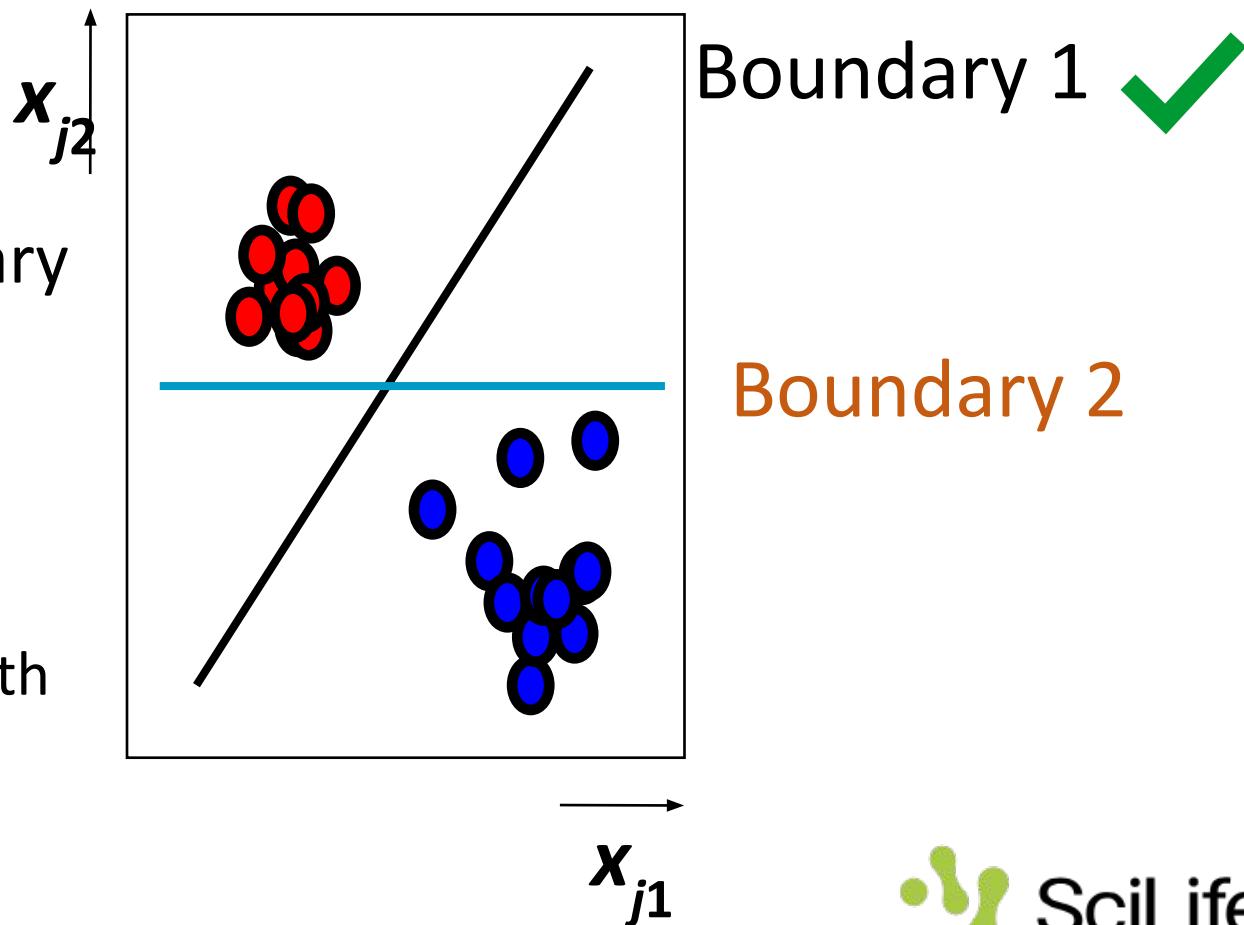
# Support Vector Machine (SVM)



Decision Boundary

# Support Vector Machine (SVM)

Which boundary
is better?



Boundary 1

Boundary 2

$x_{j2}$

$x_{j1}$

# Support Vector Machine (SVM)



Which boundary is better?

The one that maximizes the margins from both labels.

Boundary 1 ✔

Boundary 2

# 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)

Marker genes

*Prior knowledge*

# Can we automatically identify cell populations?

# Can we automatically identify cell populations?

# Benchmark paper 2019

## A comparison of automatic cell identification methods for single-cell RNA sequencing data

Tamim Abdelaal, Lieke Michielsen, Davy Cats, Dylan Hoogduin, Hailiang Mei, Marcel J. T. Reinders & Ahmed Mahfouz ✉

(Abdelaal et al. Genome Biology 2019)

# 16 existing classifiers (April 2019)

CaSTLe

scPred

LAmbDA

scmap$_{cluster}$

Moana

SingleCell
Net

CHETAH

SingleR

Garnett
SCINA

scmap$_{cell}$

DigitalCellSorter

scID

scVI

Cell-Blast

ACTINN

**RF**

CaSTLe
LAmbDA

SingleCell
Net

**NMC**

scmap$_{cluster}$

**SVM**

scPred

Moana

**Correlation**

CHETAH

SingleR

**Others**

Garnett
SCINA
DigitalCellSorter
Cell-Blast

**kNN**

scmap$_{cell}$

**LDA**

scID

**Neural networks**

scVI

ACTINN

31

# 16 existing + 6 off-the-shelf classifiers

**_RF_**

CaSTLe
LAmbDA

SingleCell
RF   Net

**_NMC_**

scmap$_{cluster}$
NMC

**_SVM_**

scPred
Moana
SVM
SVM$_{rejection}$

**_Correlation_**

CHETAH
SingleR

**_Others_**

Garnett
SCINA
DigitalCellSorter
Cell-Blast

**_kNN_**

scmap$_{cell}$
kNN

**_LDA_**

scID
LDA

**_Neural networks_**

scVI
ACTINN

# Experiment 1: intra-dataset evaluation

- Stratified 5-fold cross validation

- Performance evaluation
  - Median F1-score: $F1 = 2 \frac{precision.recall}{precision+recall}$
  - % unlabelled cells

# Most classifiers work well



Median F1-score — % Unlabeled

34

# Performance drops with deeper annotation
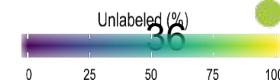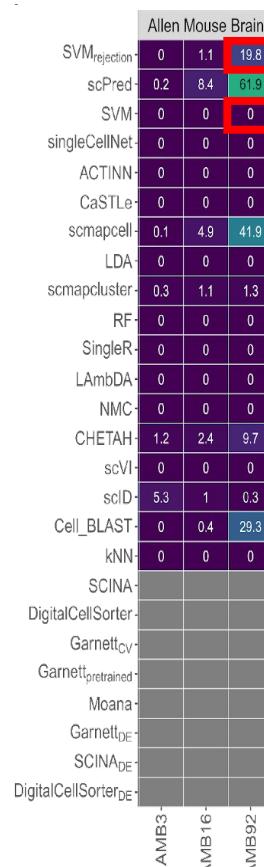


Median F1-score

% Unlabeled

# Trade-off between high performance and rejecting cells

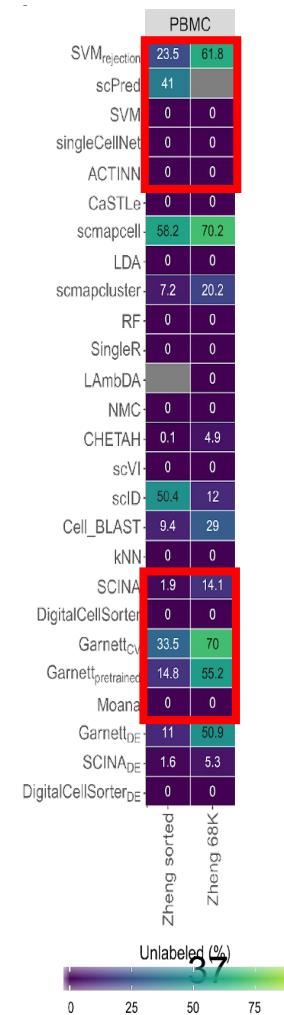## Median F1-score



## % Unlabeled

# Prior knowledge is not always beneficial

## Median F1-score
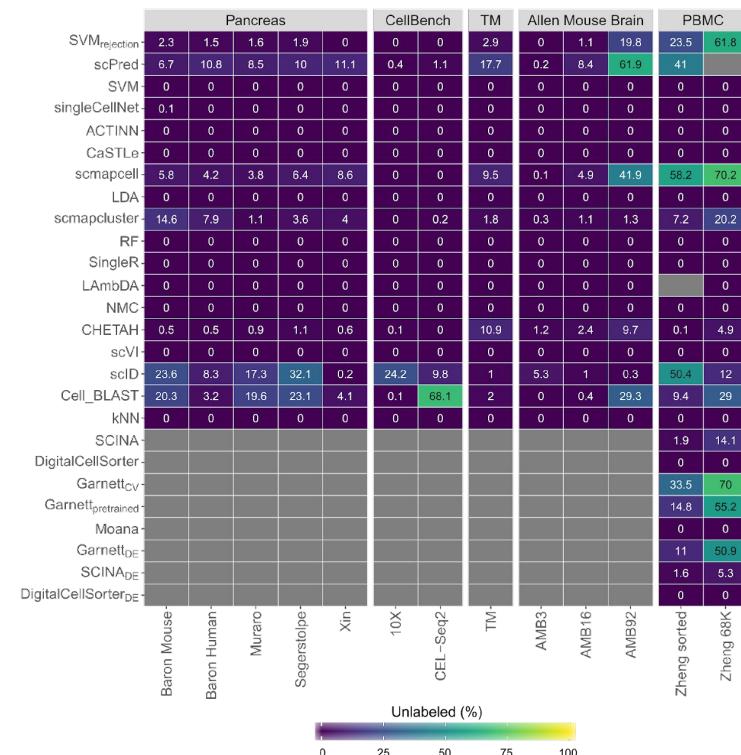


## % Unlabeled



**Lower number of classes!**

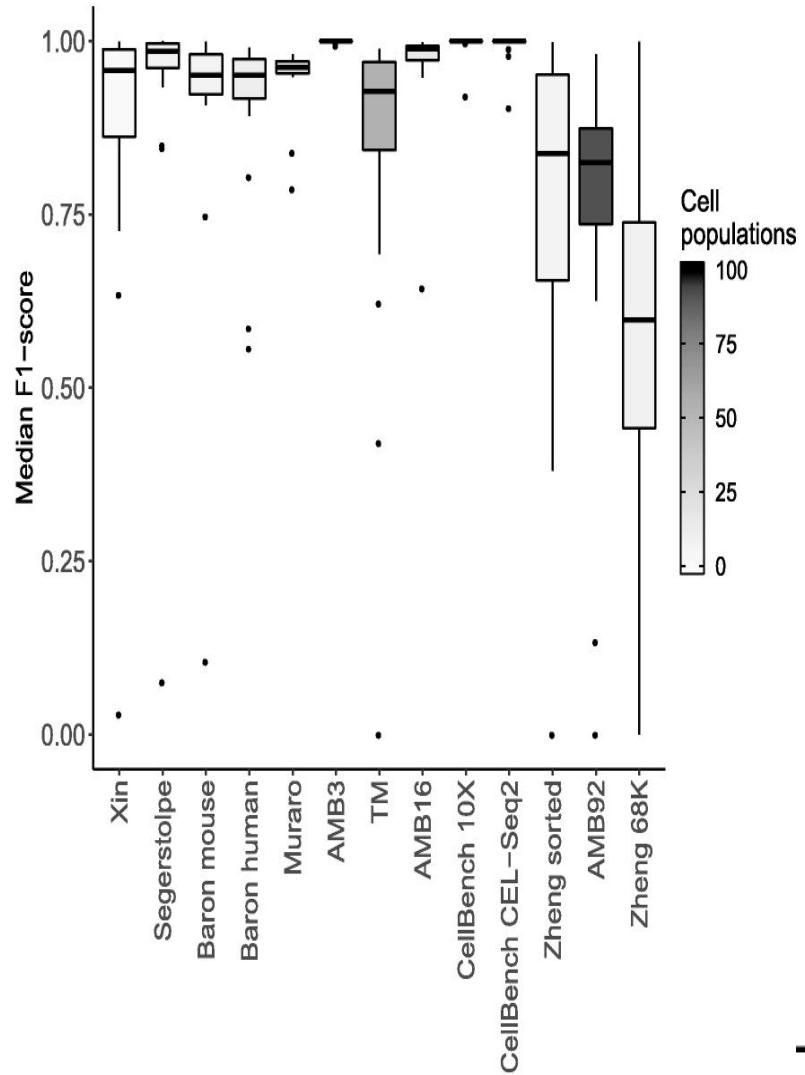# Off-the-shelf SVM outperforms dedicated single cell classifiers
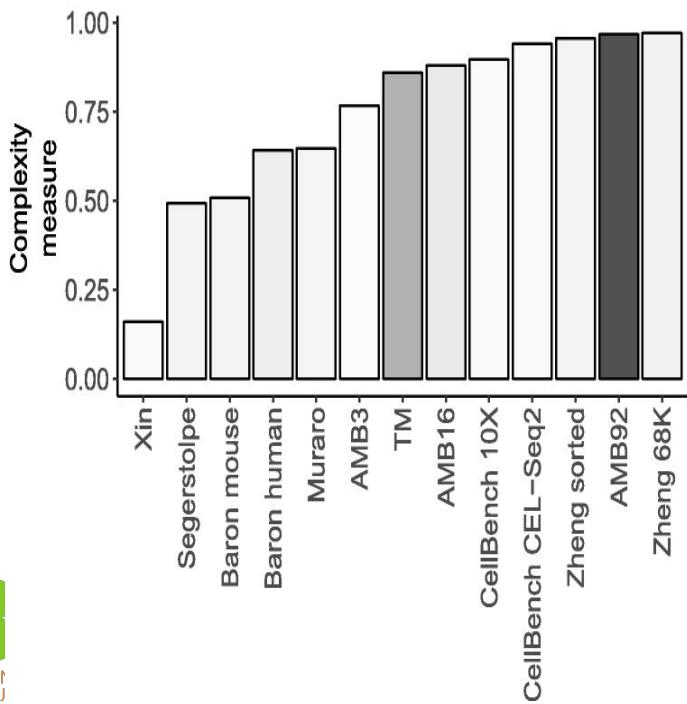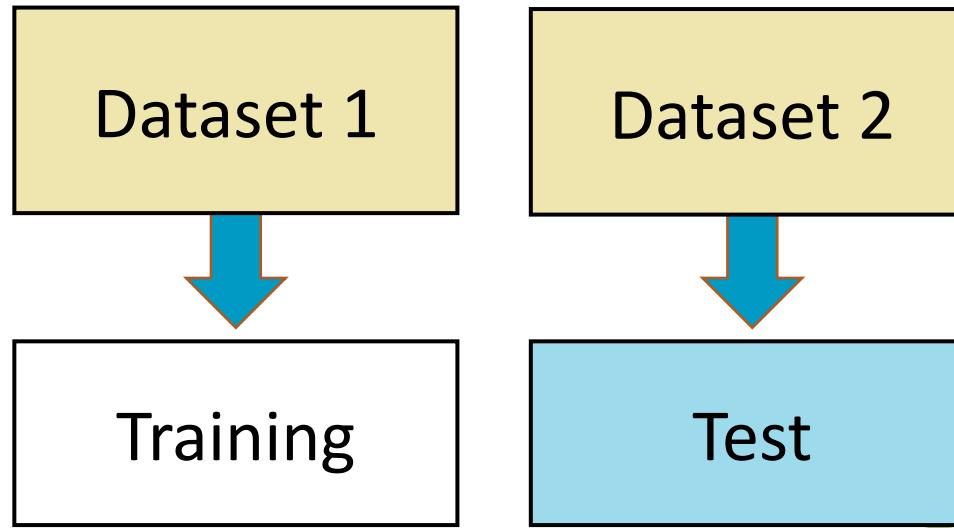


Median F1-score

% Unlabeled

# Performance depends on dataset complexity

$$\text{Complexity} = \text{mean}\left(\max_{\forall i, i \neq j} \text{corr}_{\forall i,j}\left(\text{avg}_{C_i}, \text{avg}_{C_j}\right)\right)$$
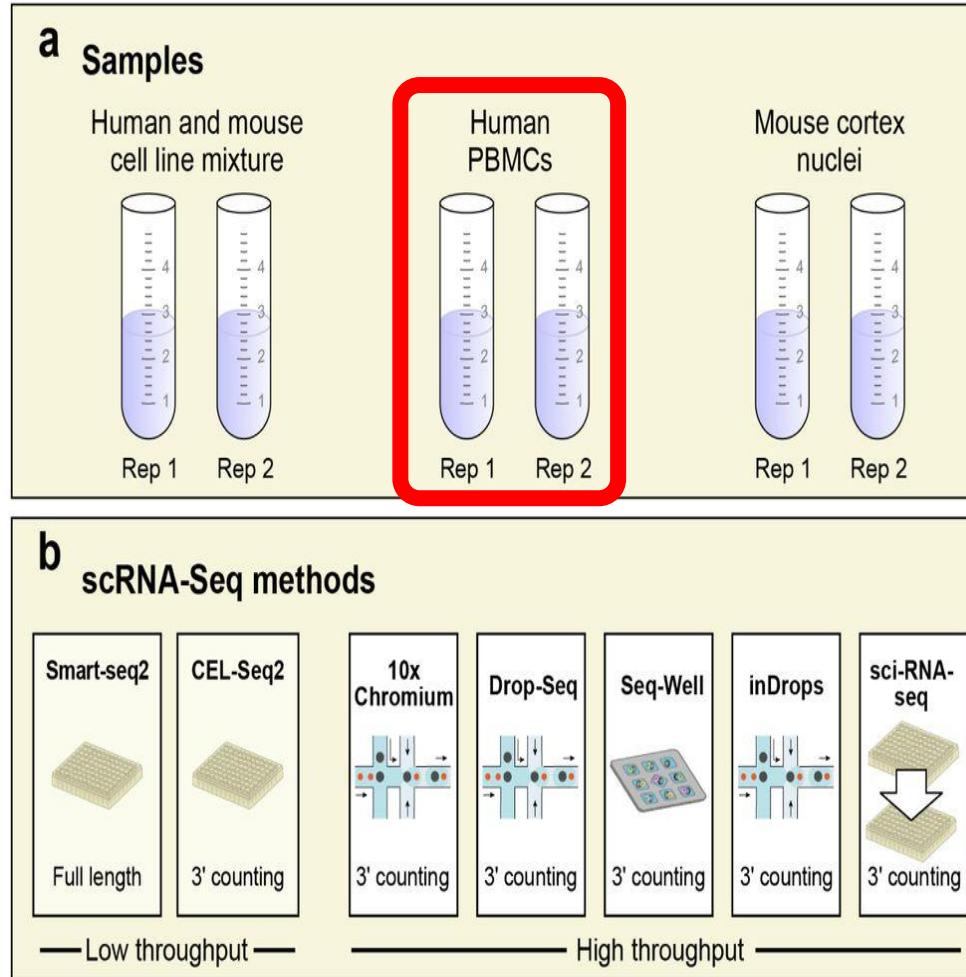
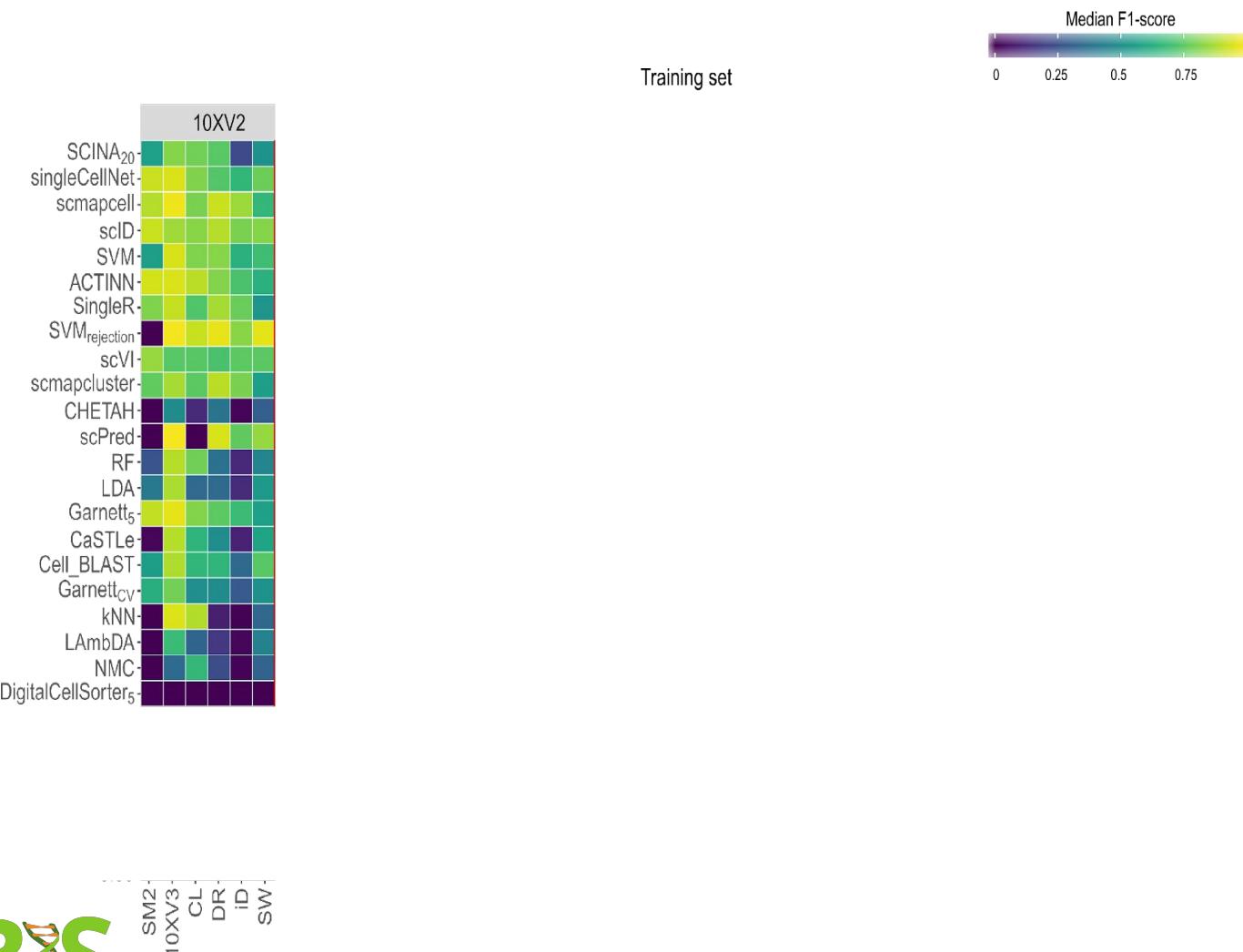# Experiment 2: inter-dataset evaluation

- Train on one dataset, evaluate on another

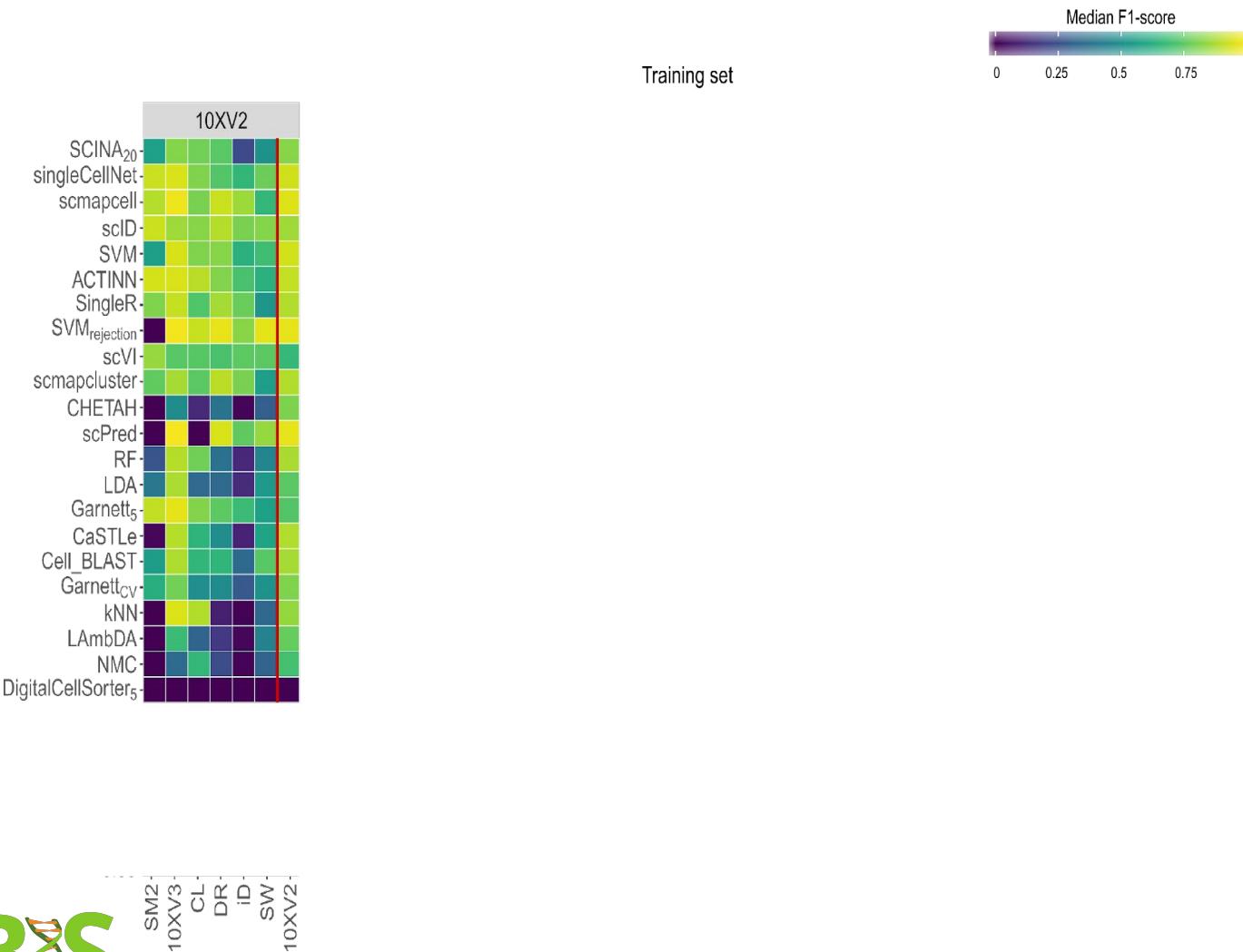- More realistic scenario

- More challenging, data is not aligned
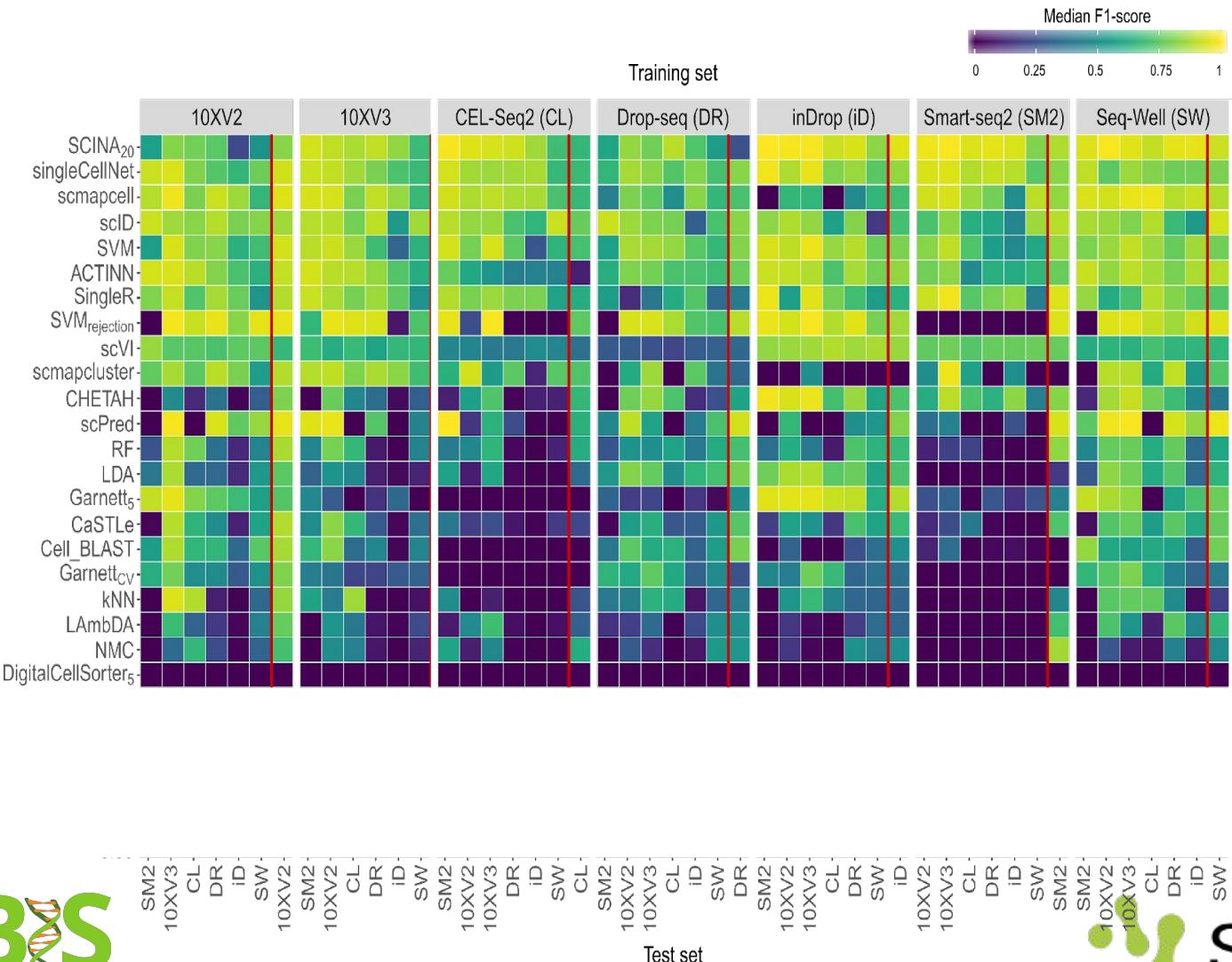
# Experiment 2: inter-dataset evaluation



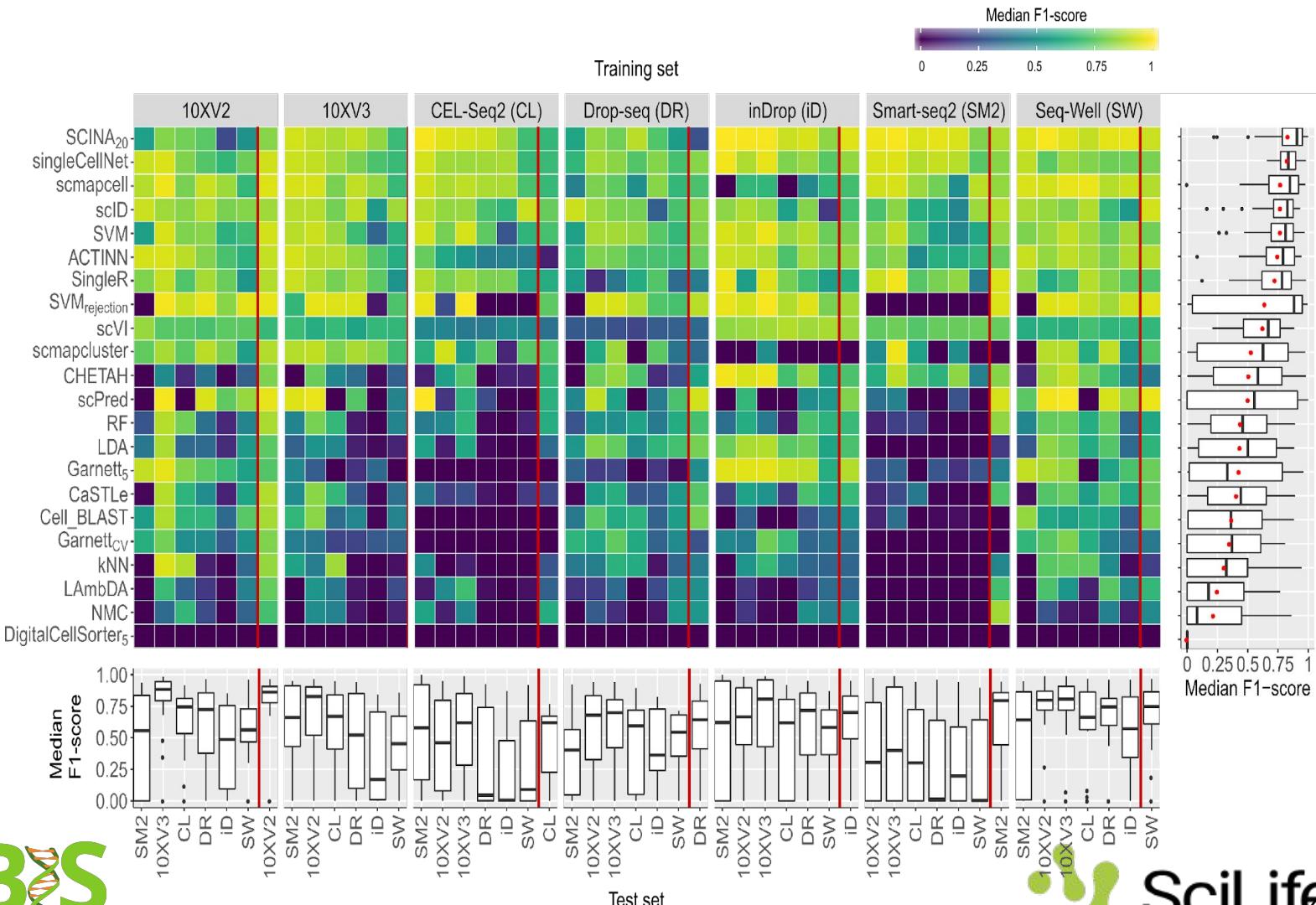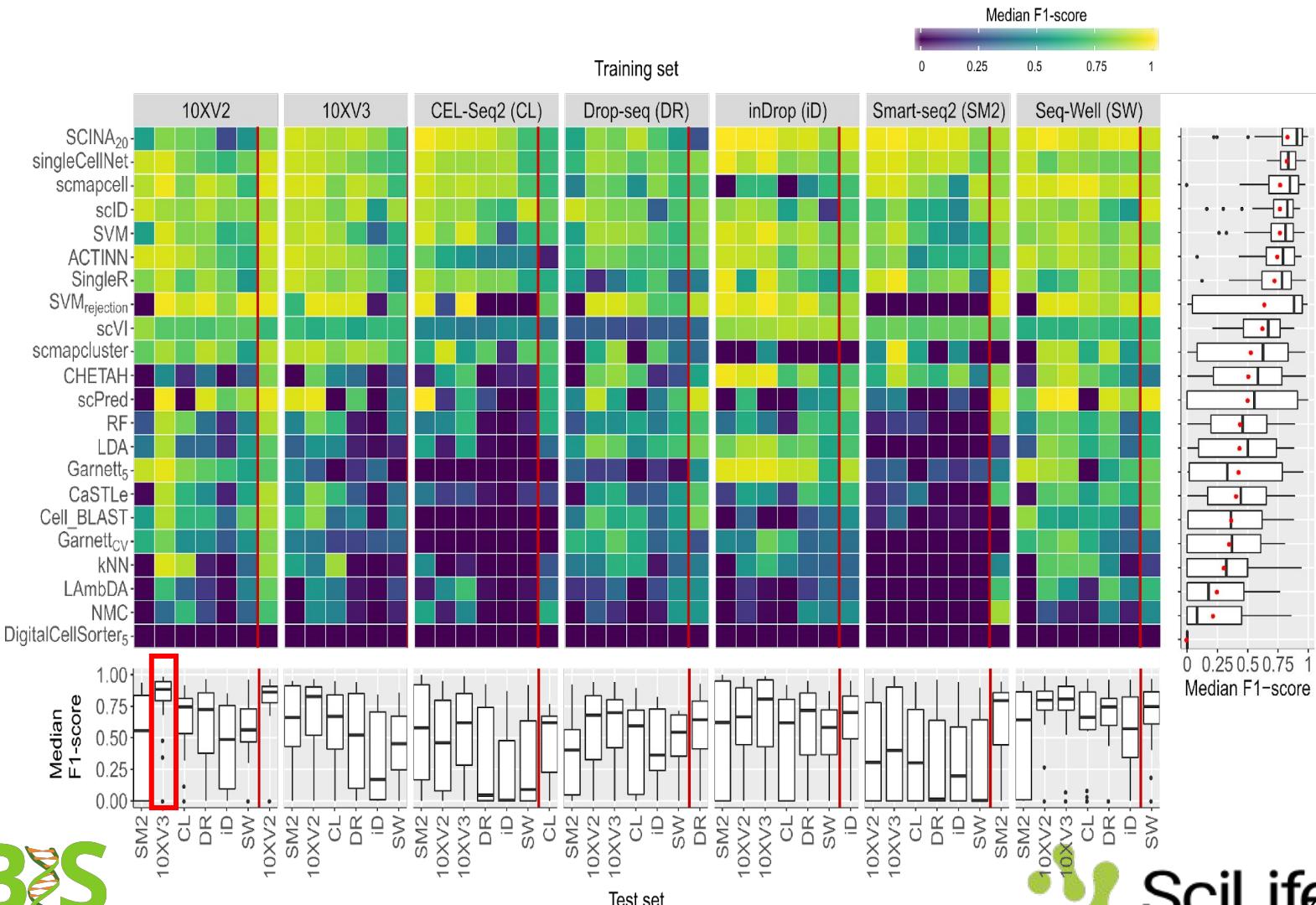Jiarui Ding et al. *Nature Biotechnology* 2020

# Prediction across protocols

# Prediction across protocols

# Prediction across protocols

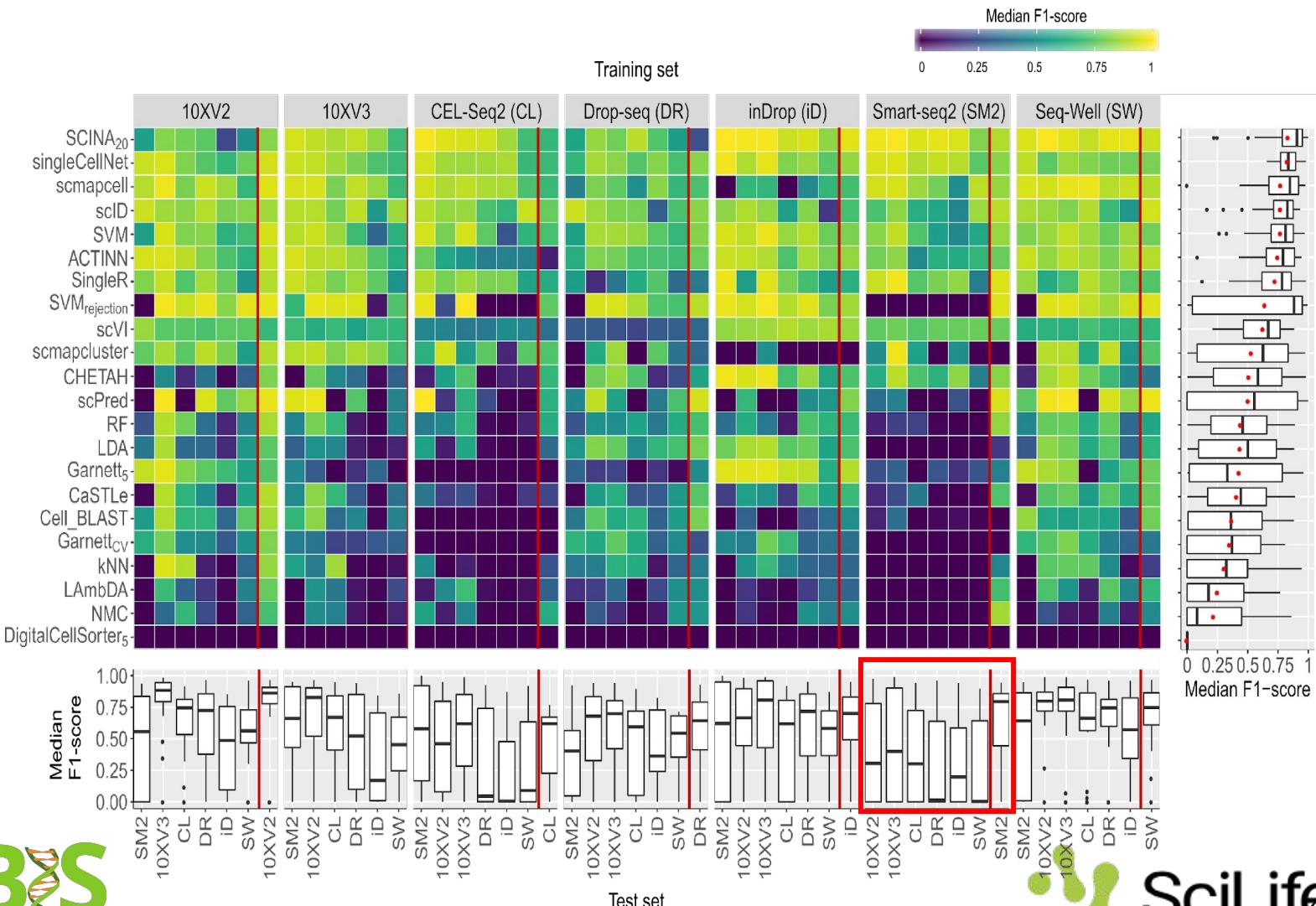# Prediction across protocols

# Prediction across protocols
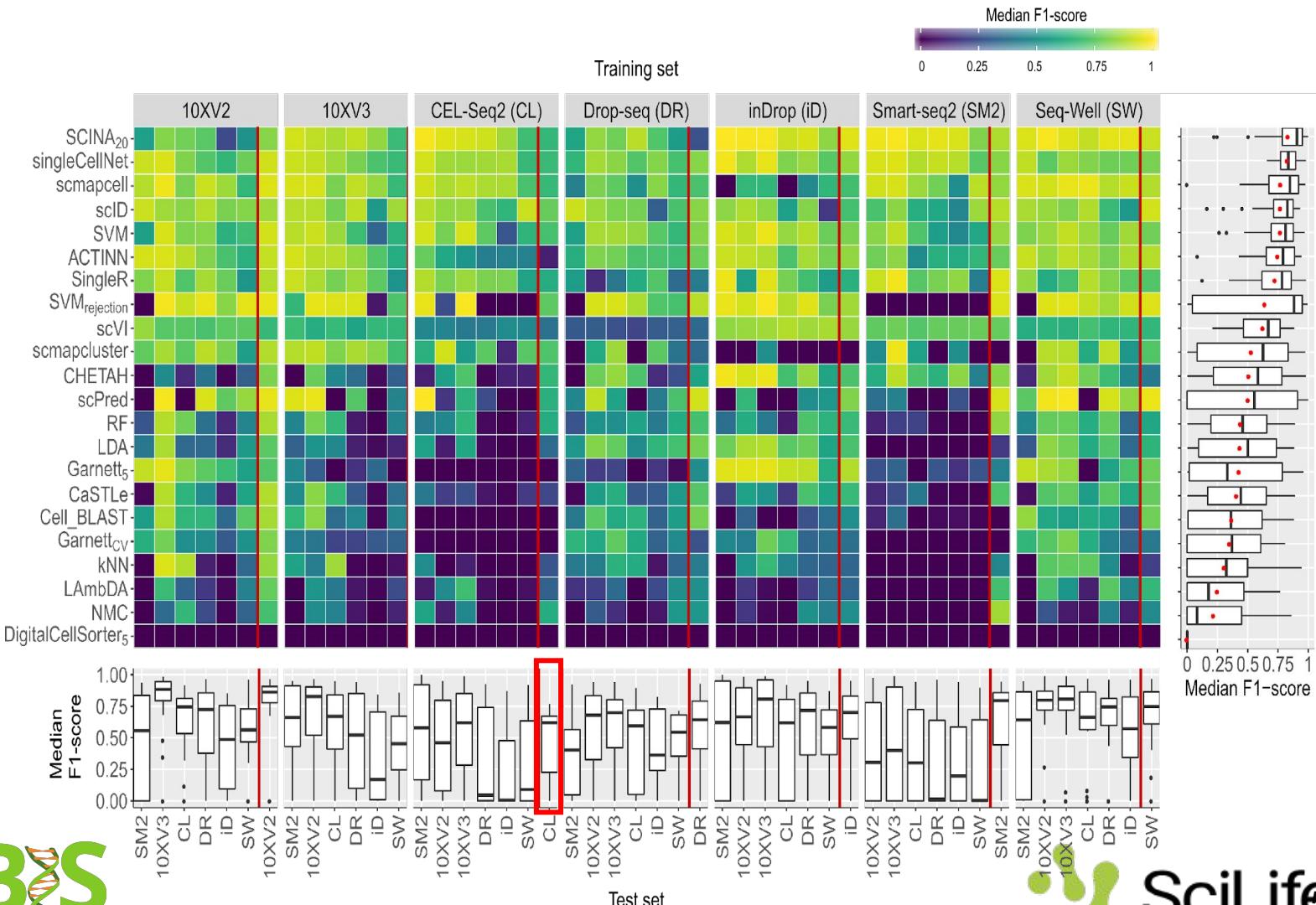
# Prediction across protocols

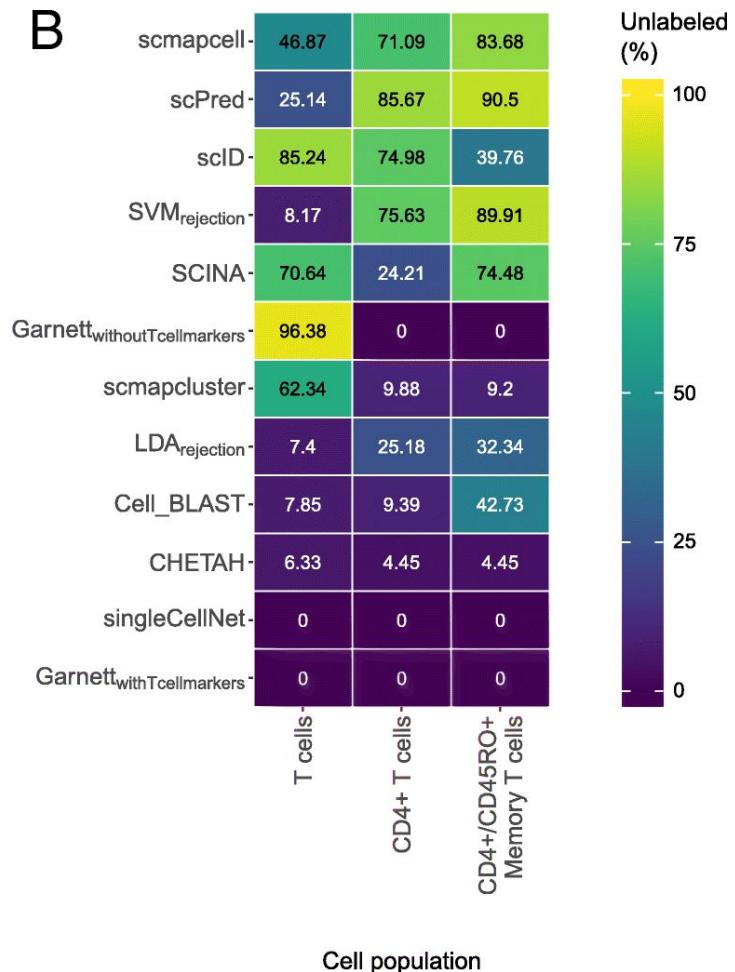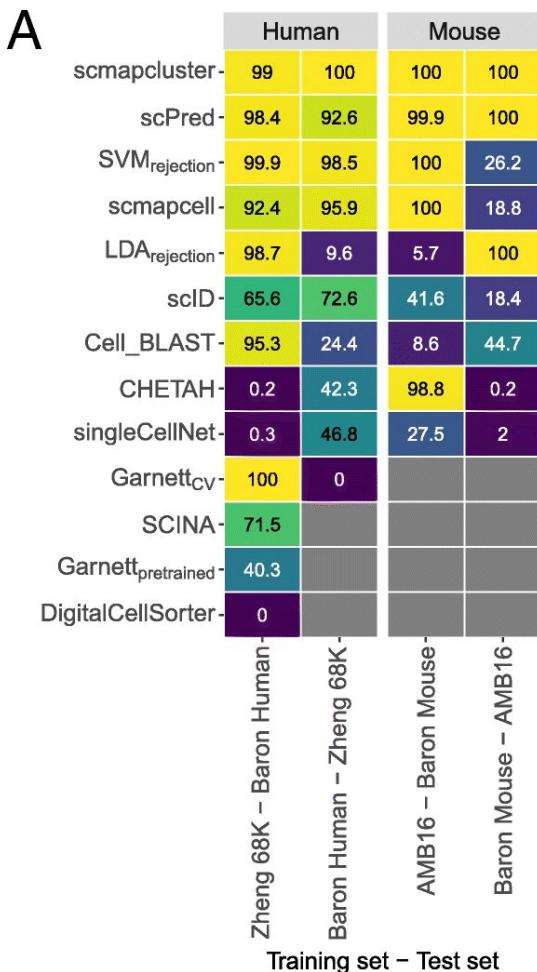# Prediction across protocols
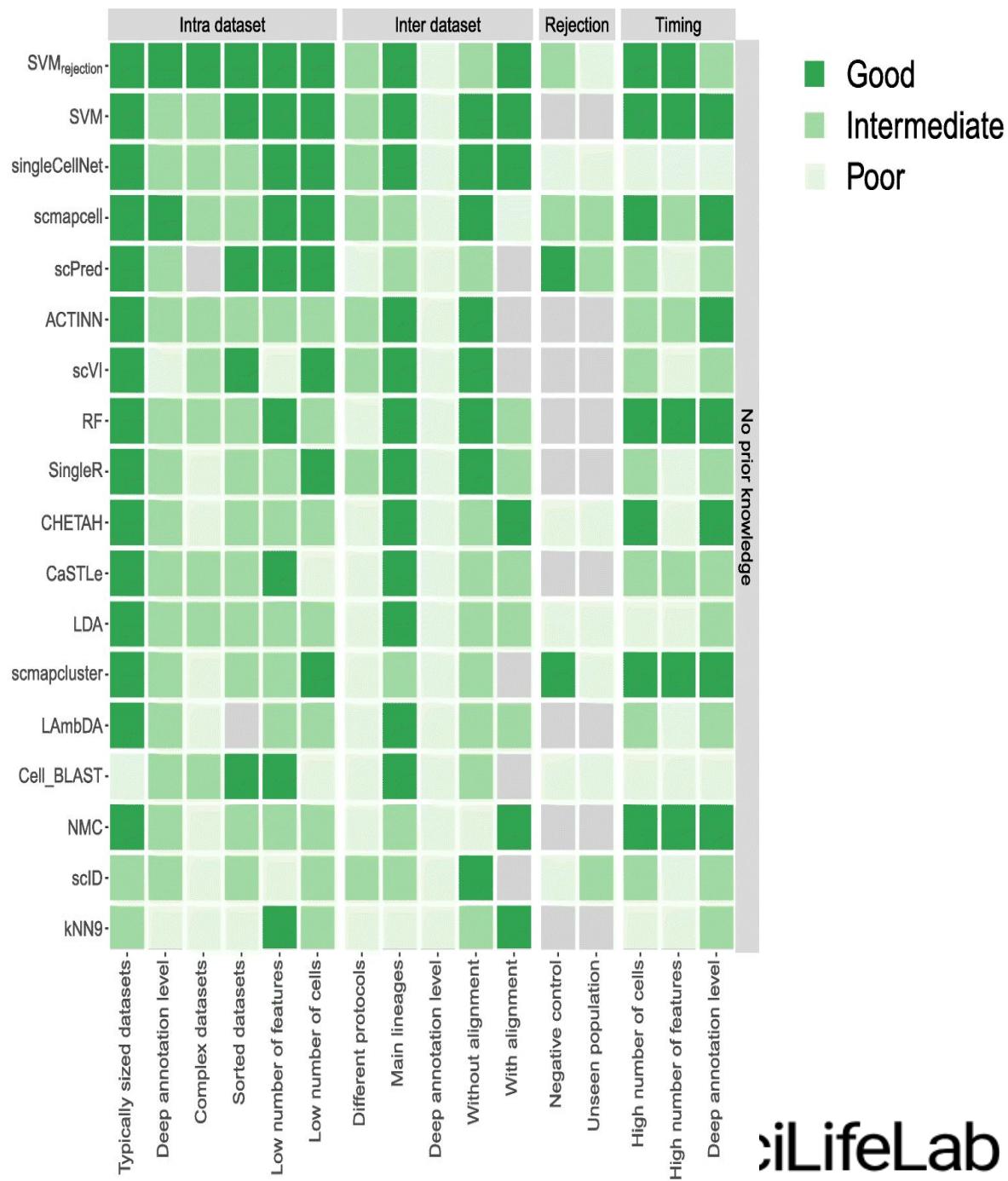
# Prediction across protocols

# Prediction across protocols

# Experiment 3: rejection evaluation

# 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/

Abdelaal* Michielsen* et al. *Genome*

# Benchmark paper 2021
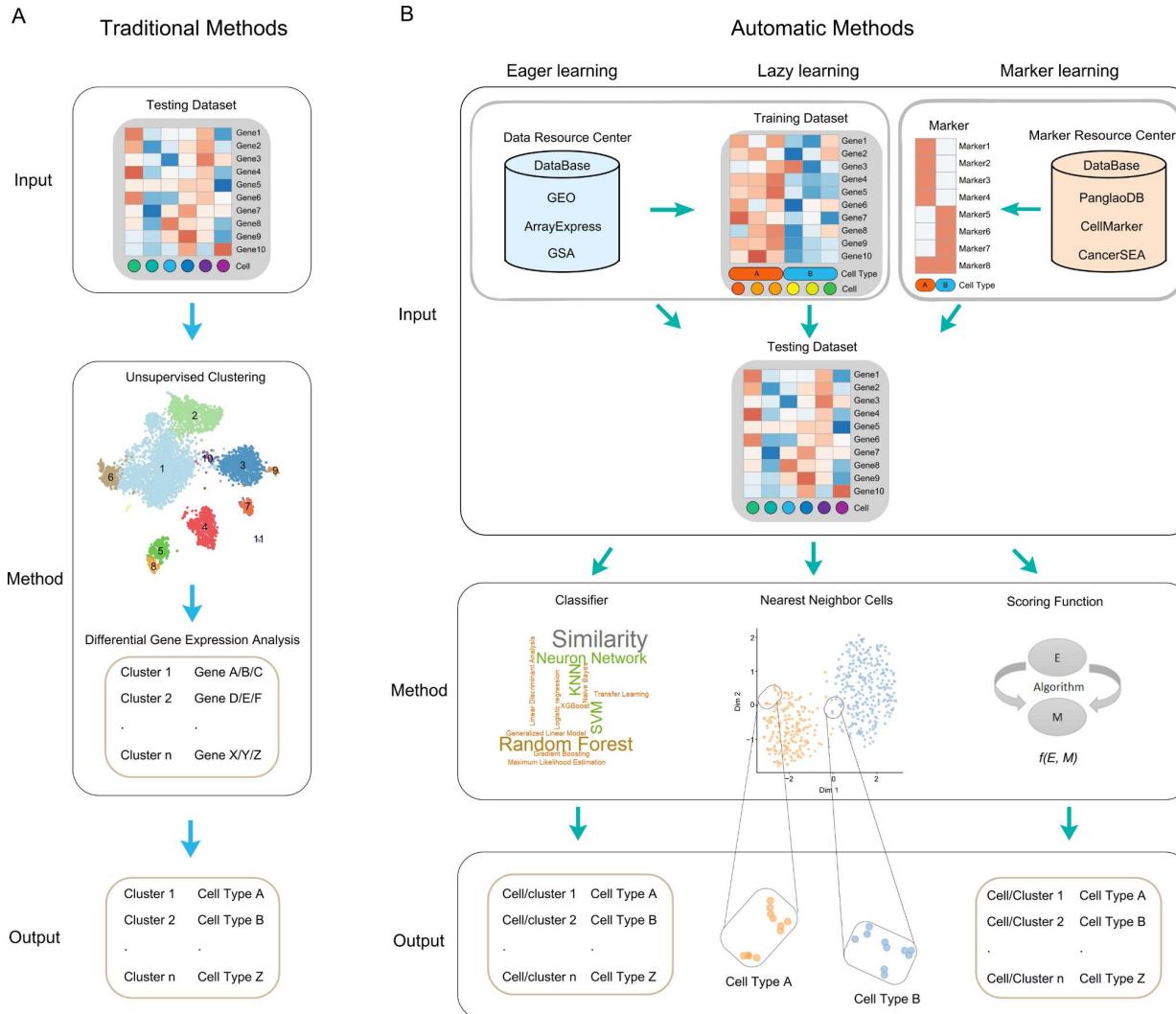


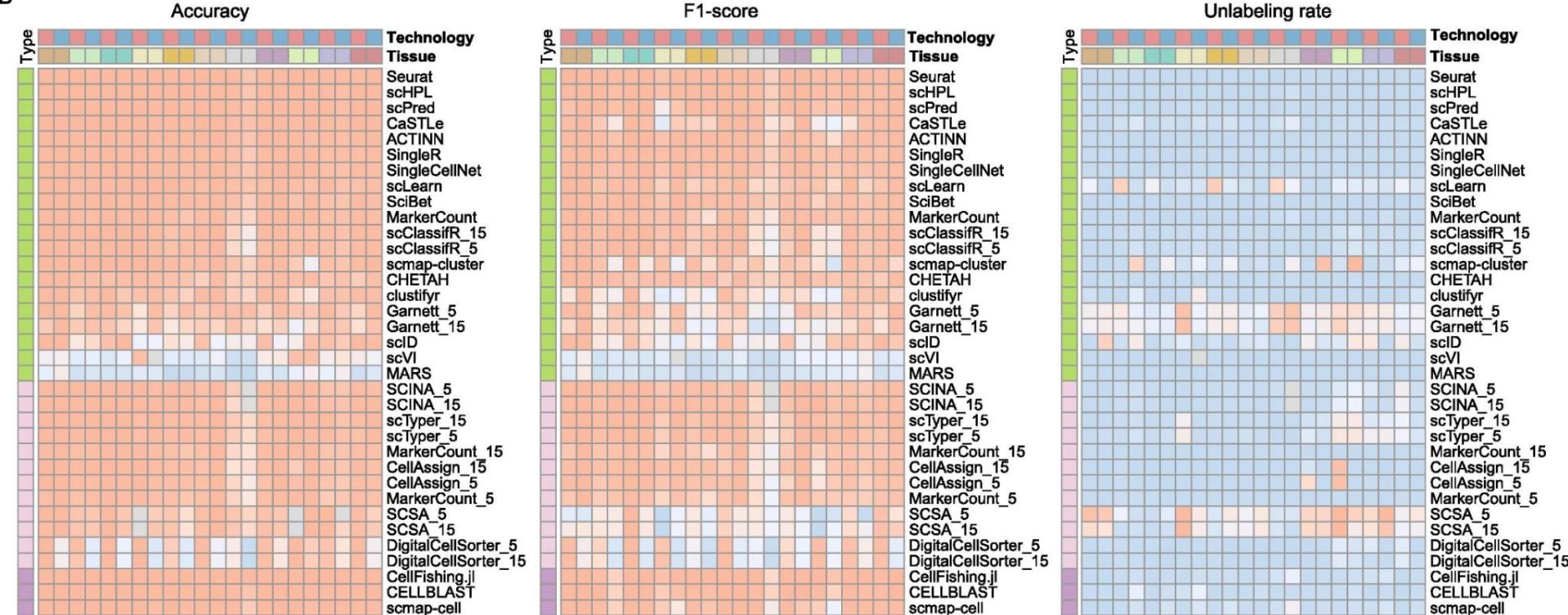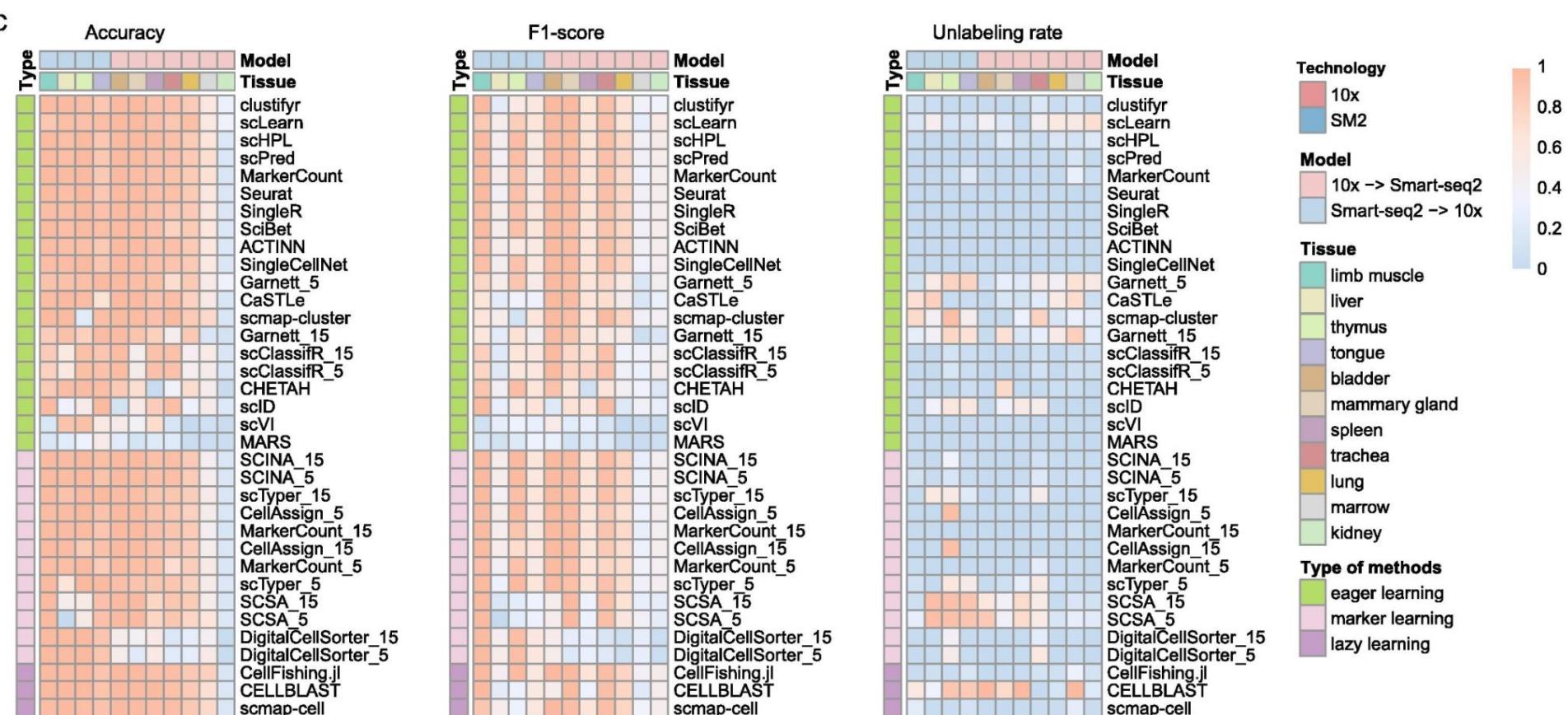(Xie et al. Comp. Struct. Biotech J. 2021)

# Table with all the methods

https://www.csbj.org/action/showFullTableHTML?isHtml=true&tableId=t0005&pii=S2001-0370%2821%29004 49-9

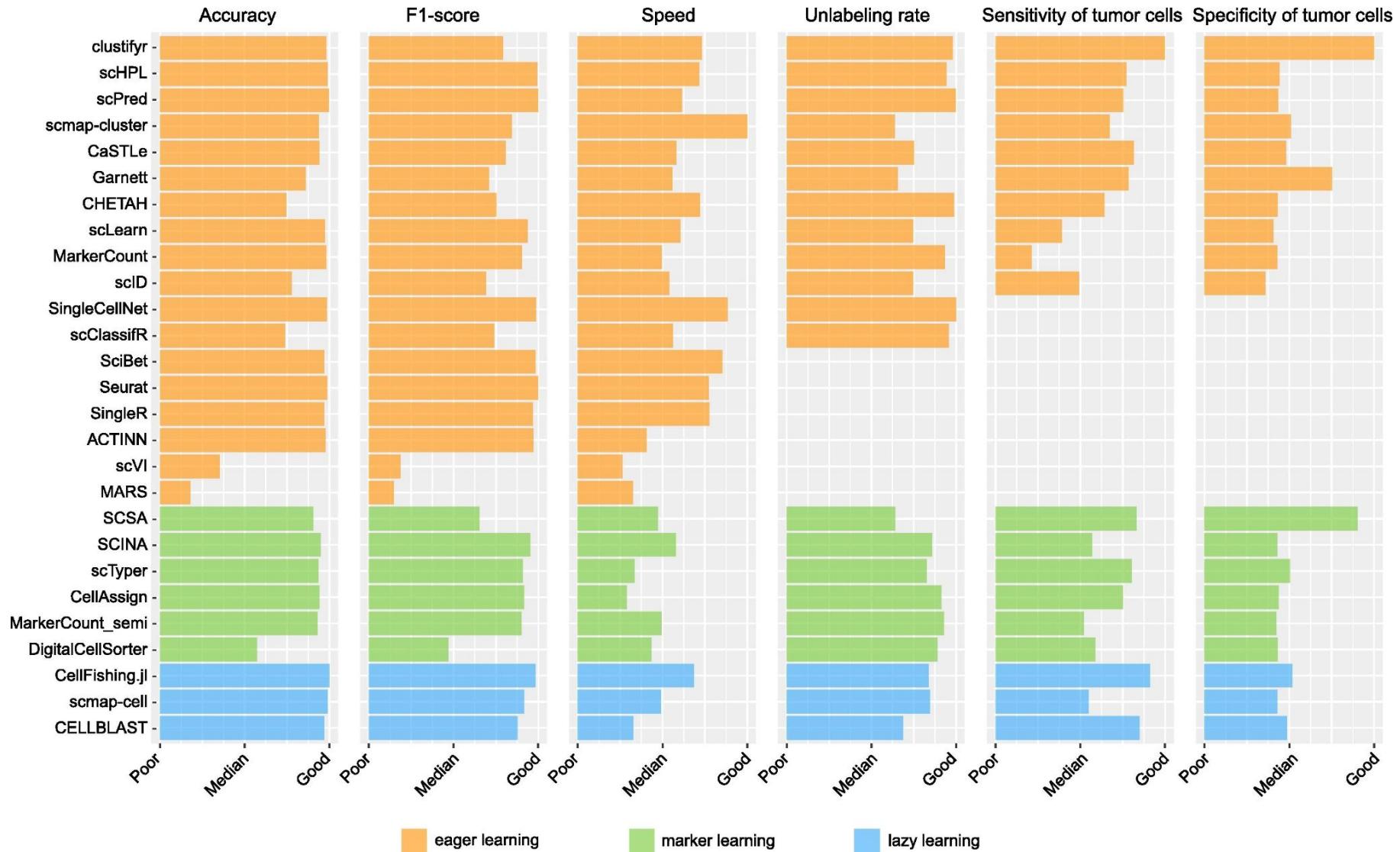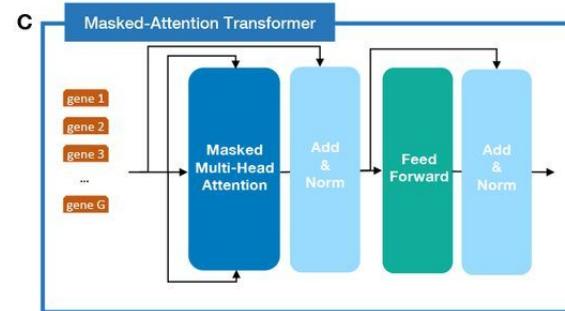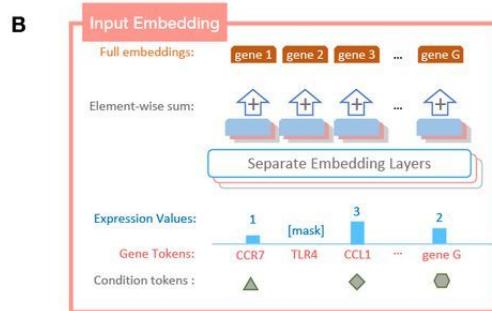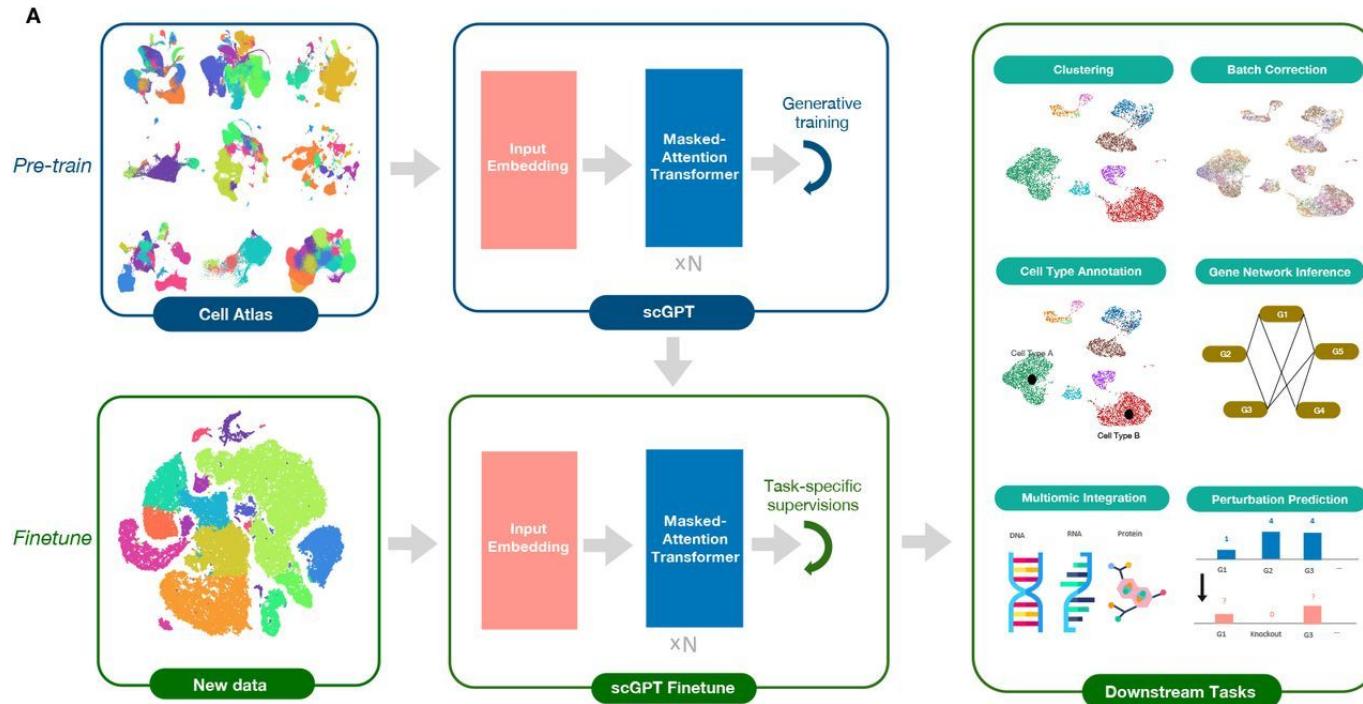# Within dataset training/testing with cross-validation
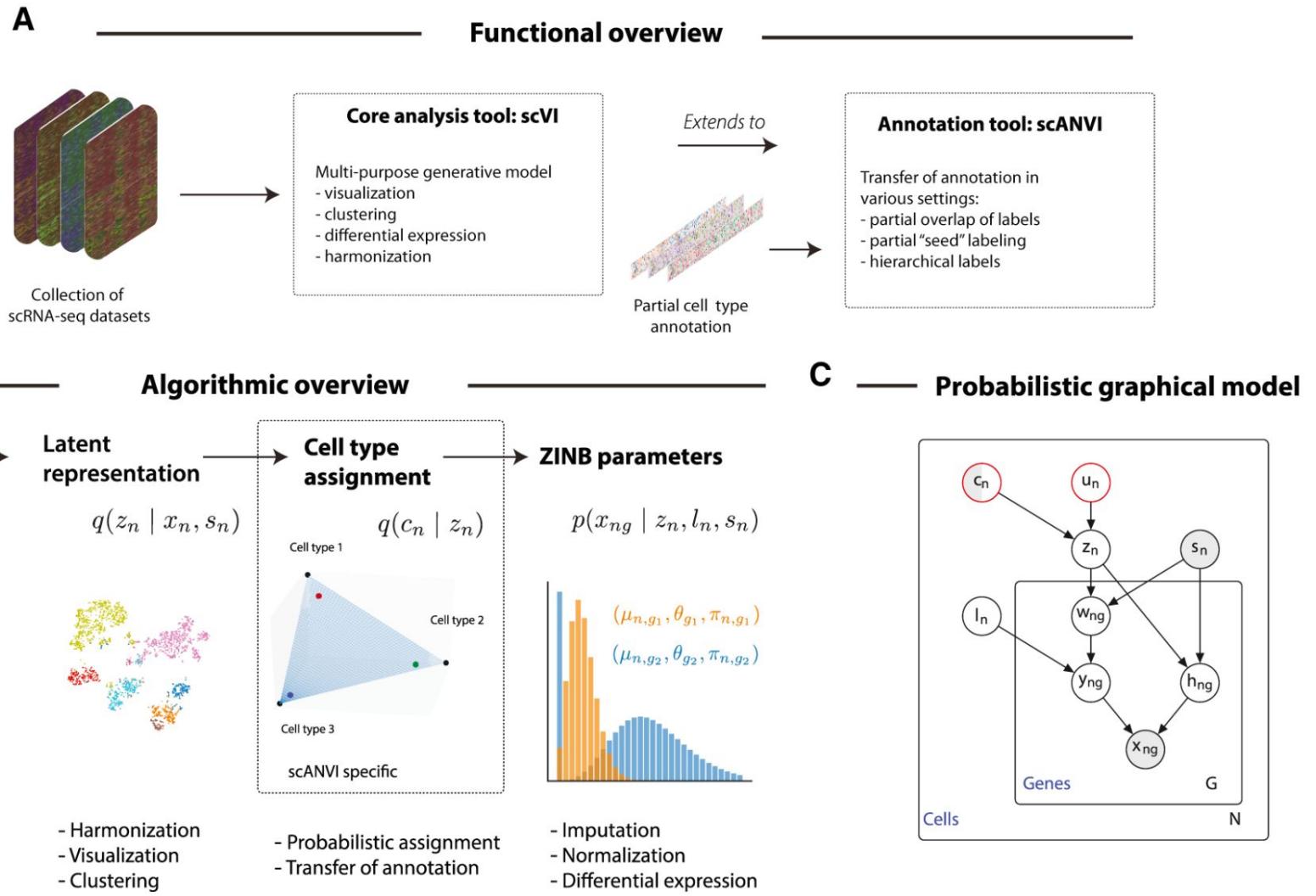
# Across technologies

# Summary

# Generative learning is the next big thing? scGPT

# Generative learning is the next big thing? scANVI

# Some useful resources

- Azimuth - Seurat label transfer to reference sets
  - https://azimuth.hubmapconsortium.org/
  - online or R package
- DISCO - CellMapper to several tissues
  - https://www.immunesinglecell.org/
- Celltypist - Regularised linear models with Stochastic Gradient Descent
  - https://www.celltypist.org/
  - online or python package

# Summary

- Cell identification is moving from unsupervised (clustering/visualization) to supervised (classification) learning
- Check what reference you are using!
  - The more similar reference is to your data - the better the prediction.
  - Same technology matters
  - Do you trust their celltype annotations?
- Atlases do not contain all tissues/celltype and especially not all disease states of cells.
- Also look at DGE and known markers and check that predictions makes sense