Celltype Prediction

Fariba Roshanzamir 03-Apr-2025

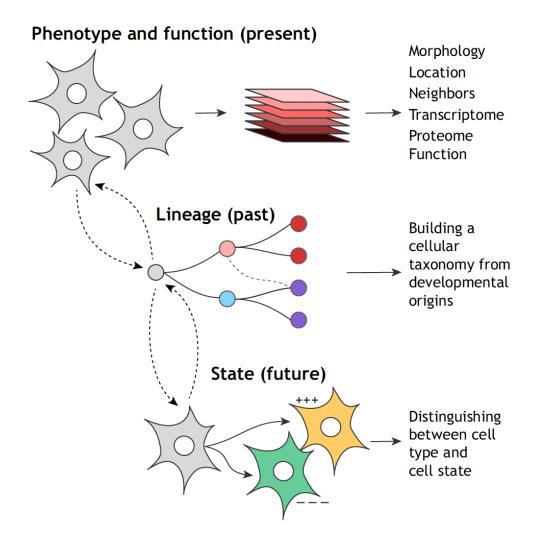
fariba.roshanzamir@scilifelab.se

Åsa Björklund (NBIS, scilifelab) Ahmed Mahfouz (Leiden University Medical Center, TU Delft)





Cell identity



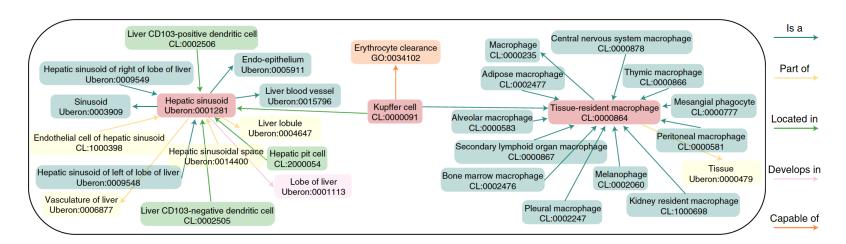




Celltype ontologies

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

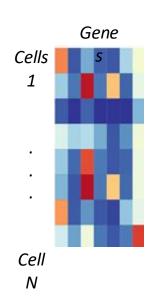
Including HuBMAP, Human Cell Atlas (HCA), cellxgene, Single Cell Expression Atlas, BRAIN Initiative Cell Census Network (BICCN), ArrayExpress, The Cell Image Library, ENCODE, and FANTOM5,

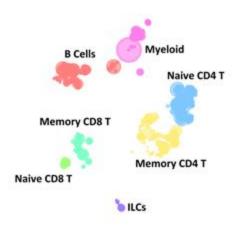






How can we identify cell populations?

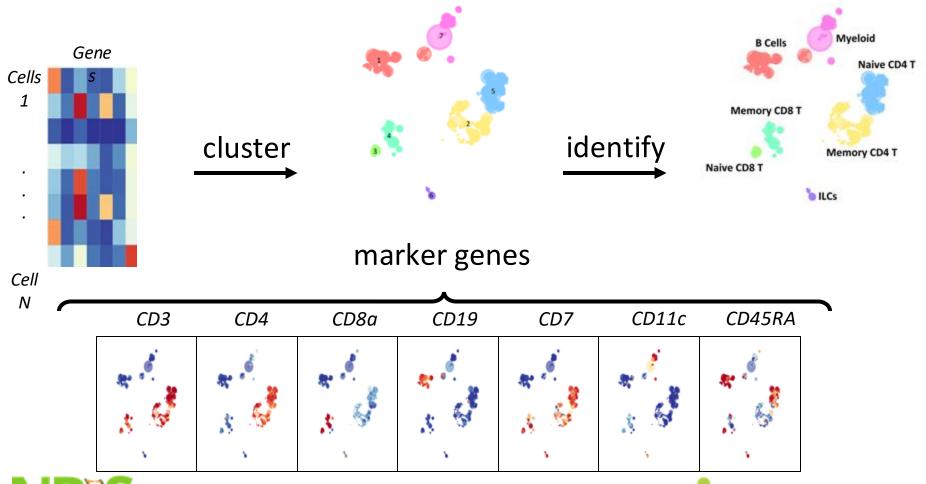








How can we identify cell populations?





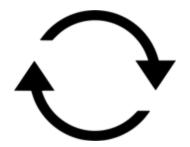


Unsupervised celltype identification is problematic

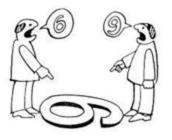
Time consuming



Not reproducible

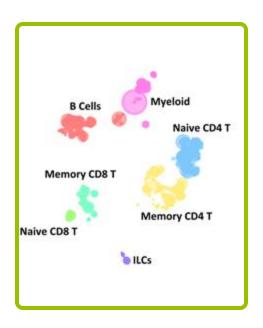


Subjective



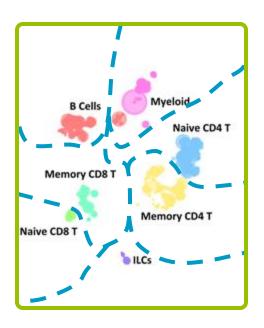










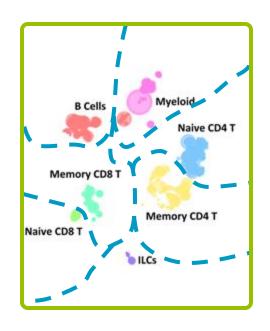






Clustering

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



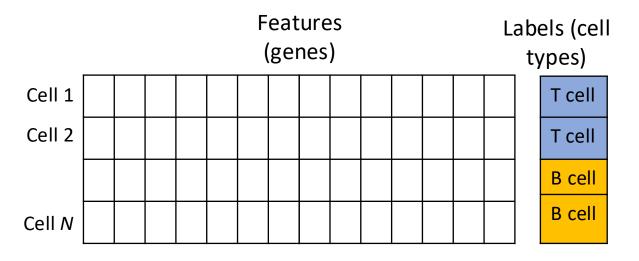
Classification

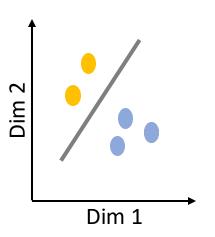
- Supervised learning
- Prior information available about different groups
- Classifiers find descriptions of decision boundaries

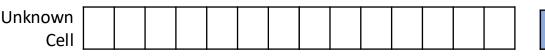




Classification













Classifier training

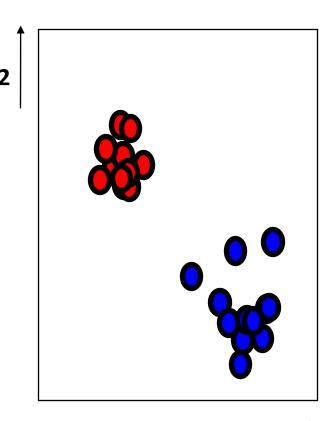
- Dataset: for *j* th cell:
 - gene expressions \mathbf{x}_{j}
 - class label: $y_j \in \{1=T, -1=B\}$
- Classifier:

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

• Errors: $E = \text{sum}(E_j)$ $E_j = \begin{cases} 1 & \text{if } \hat{y}_j \neq y_j \\ 0 & \text{if } \hat{y}_j = y_j \end{cases}$











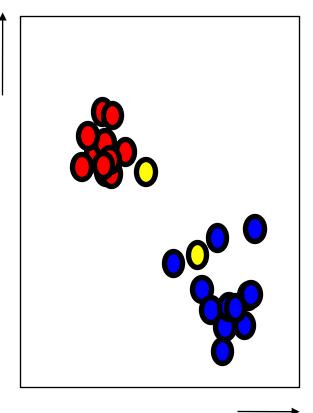
Instance Based Learning (Lazy Classification)

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



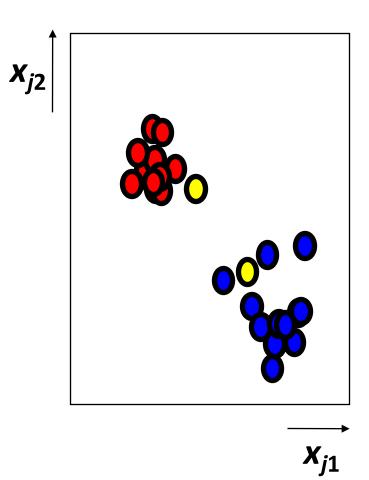






Nearest Neighbor (k-NN)

- 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?

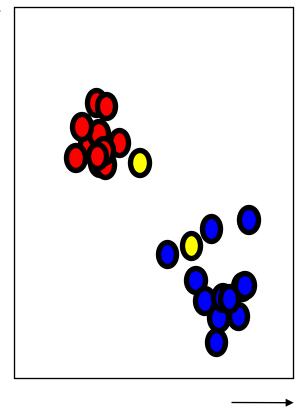






Nearest Neighbor (k-NN)

- 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

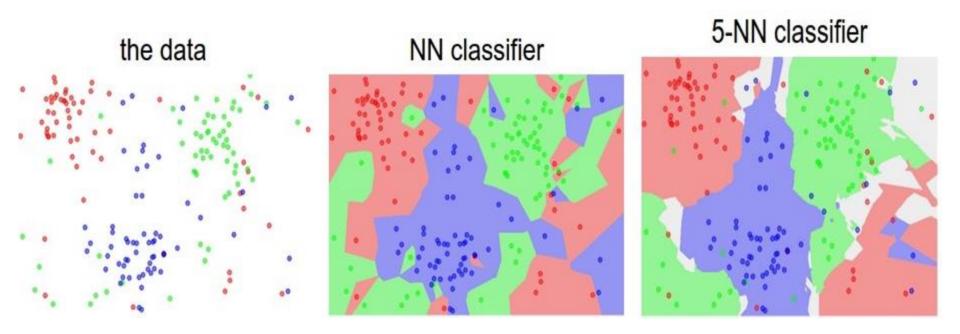








Effect of k





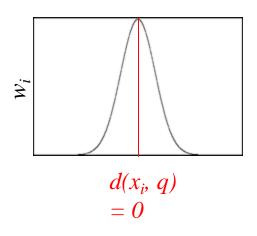


Weighted Nearest Neighbor (kernel regression)

- Distance metric:
 - Euclidean
- How many neighbors to look at?
 - All of them!
- Weighting function:

$$w_i = \exp(-\frac{d(x_i,q)^2}{K_w})$$

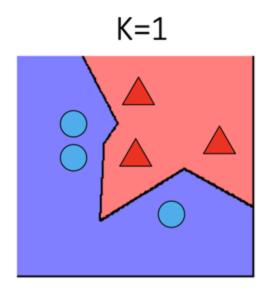
- Nearby points to a query q are weighted more strongly. K_w: kernel width
- How to fit with the local points?
 - Predict the weighted average $\frac{\sum_{i} w_{i} y_{i}}{\sum_{i} w_{i}}$

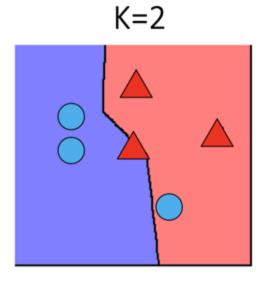


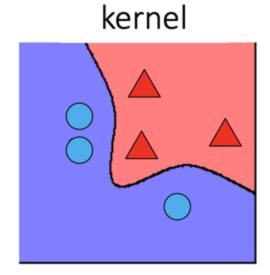




Comparison: K=1, K=2, kernel



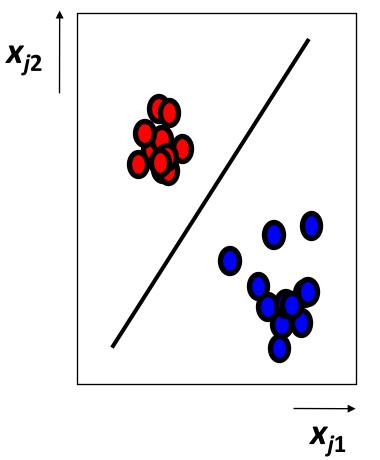








Support Vector Machine (SVM)



Decision Boundary

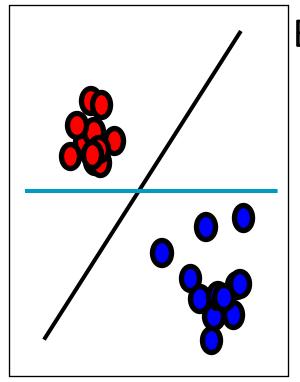




Support Vector Machine (SVM)

 X_{j2}

Which boundary is better?



Boundary 1

Boundary 2

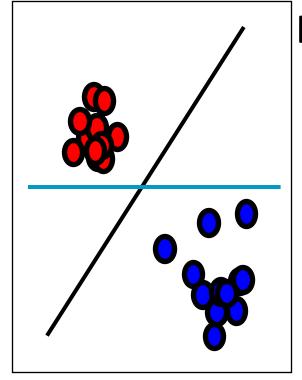




Support Vector Machine (SVM)

Which boundary is better?

The one that maximizes the margins from both labels.



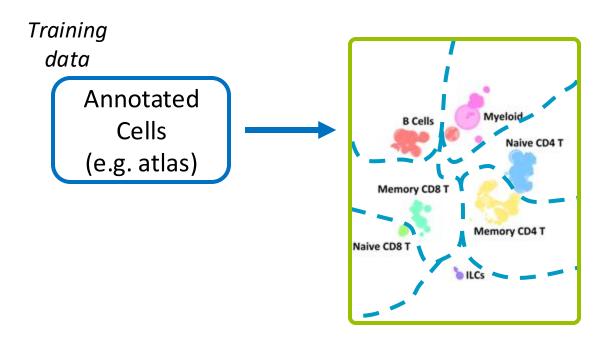
Boundary 1 🗸



Boundary 2

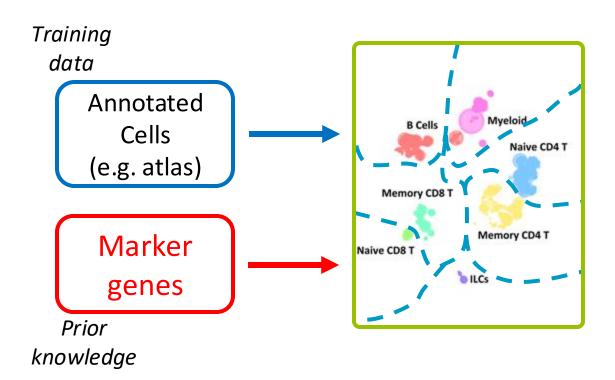






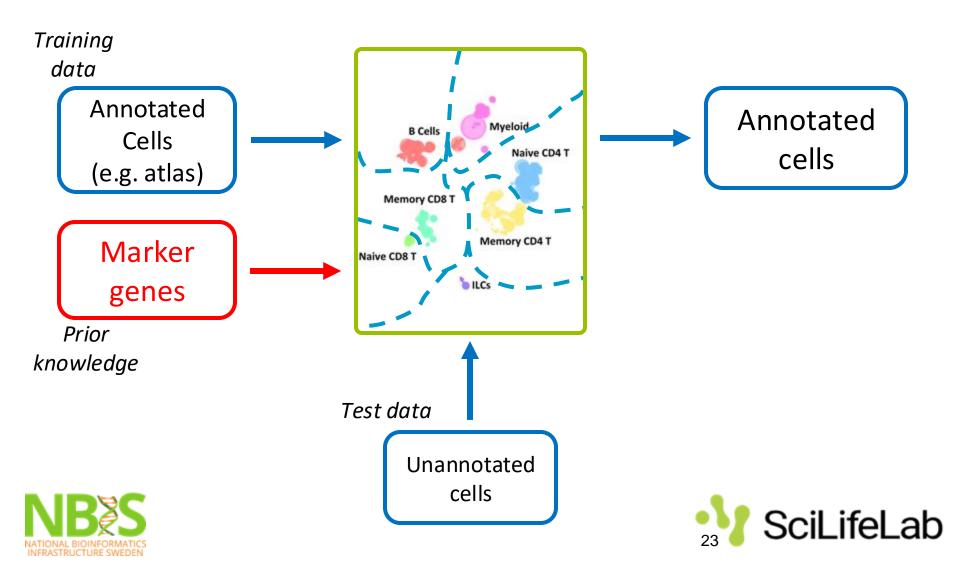


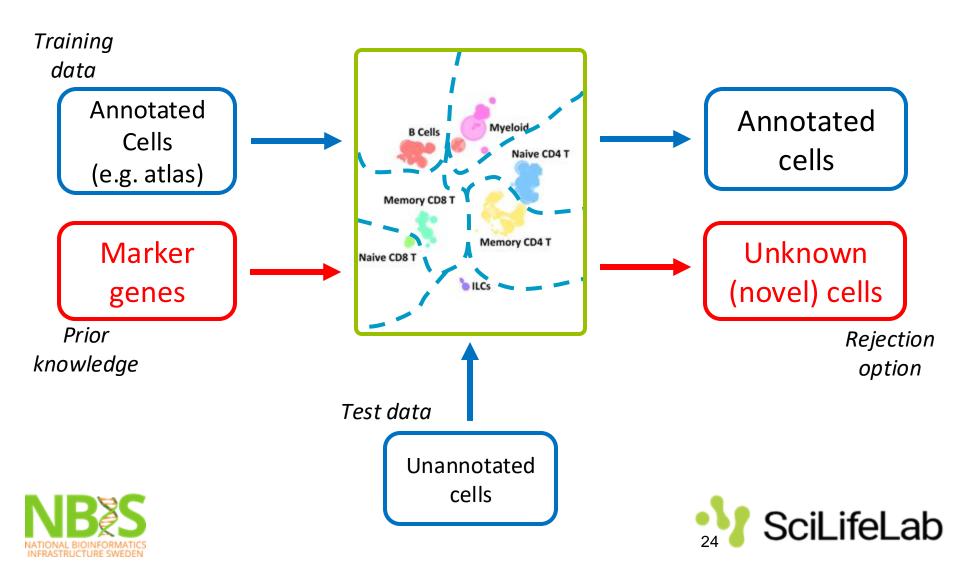












Benchmark paper 2019

Research Open access Published: 09 September 2019

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

<u>Tamim Abdelaal, Lieke Michielsen, Davy Cats, Dylan Hoogduin, Hailiang Mei, Marcel J. T. Reinders</u> & Ahmed Mahfouz □

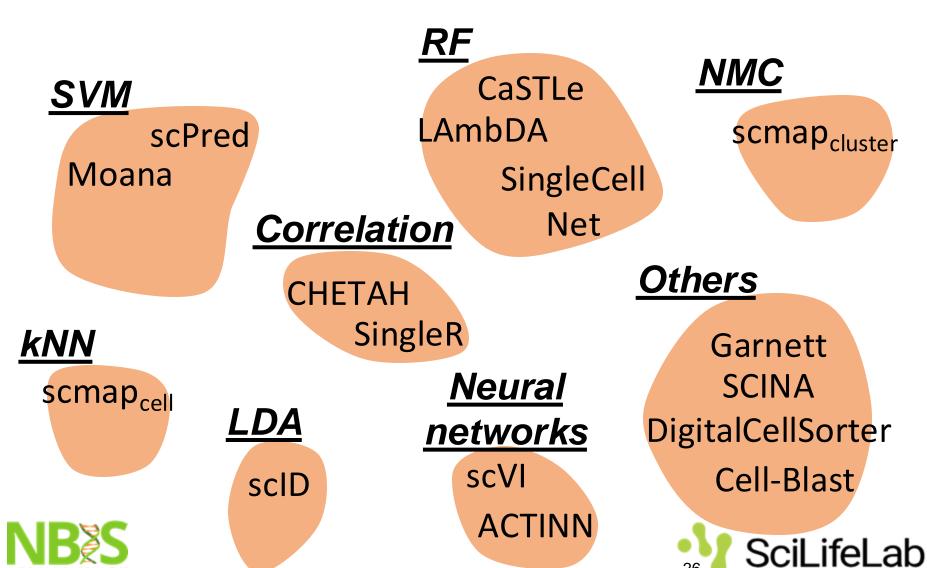
Genome Biology 20, Article number: 194 (2019) | Cite this article

COL Accesso | 277 Citations | 76 Altmotric | Matrice

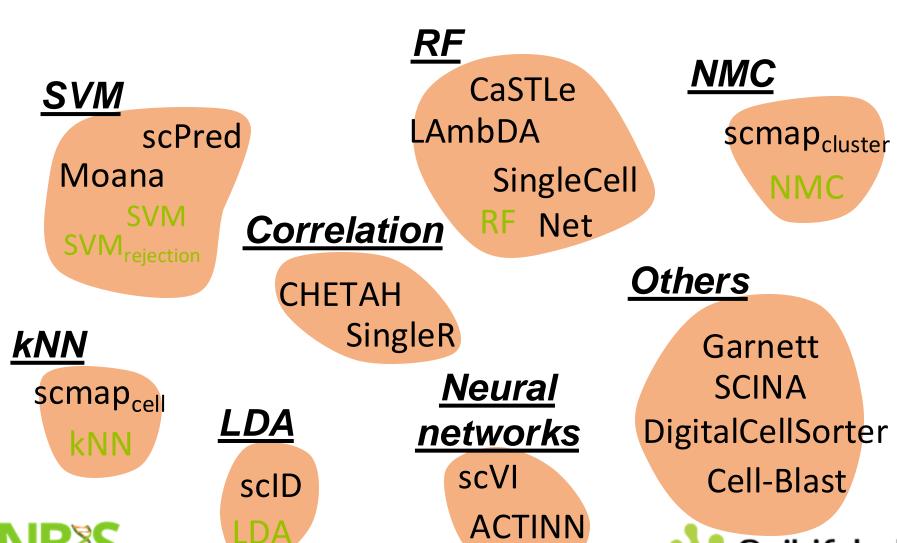




16 existing classifiers (April 2019)



16 existing + 6 off-the-shelf classifiers



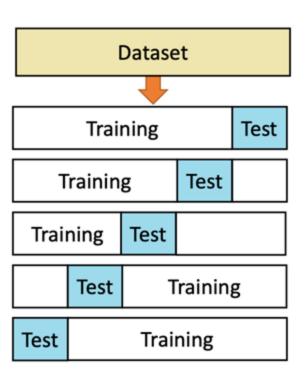
SciLifeLab

Experiment 1: intra-dataset evaluation

Stratified 5-fold cross validation

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

It ranges from 0 to 1, with 1 indicating perfect precision and recall

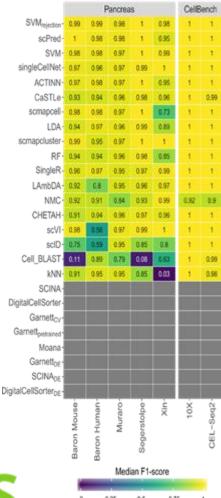


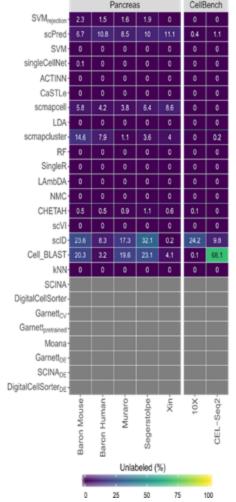




Most classifiers work well

Median F1-score



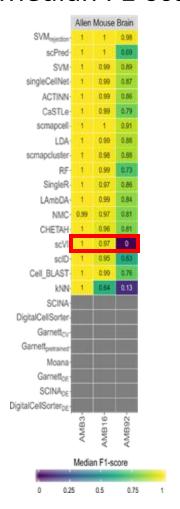


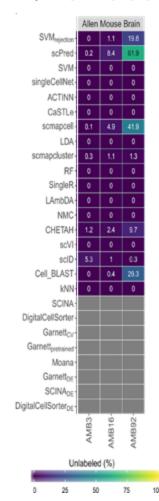




Performance drops with deeper annotation

Median F1-score



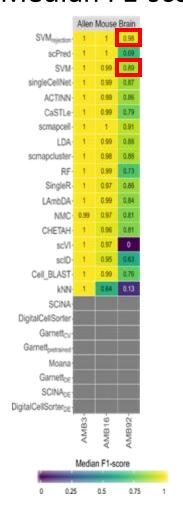


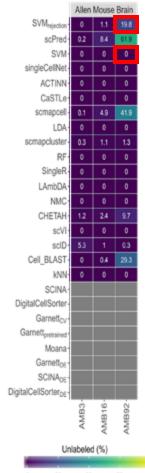




Trade-off between high performance and rejecting cells

Median F1-score



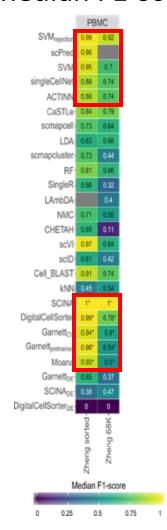




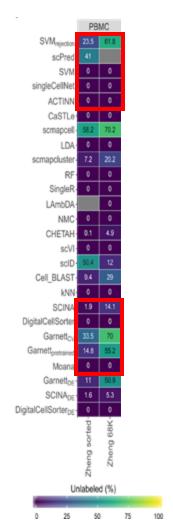


Prior knowledge is not always beneficial

Median F1-score



% Unlabeled





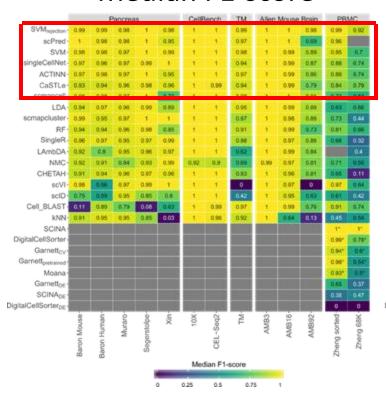
Lower number of classes!

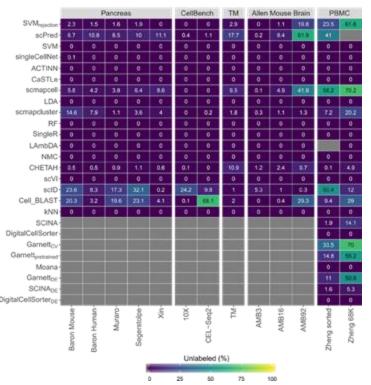




Off-the-shelf SVM outperforms dedicated single cell classifiers

Median F1-score



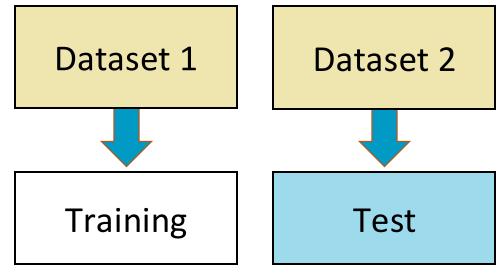






Experiment 2: inter-dataset evaluation

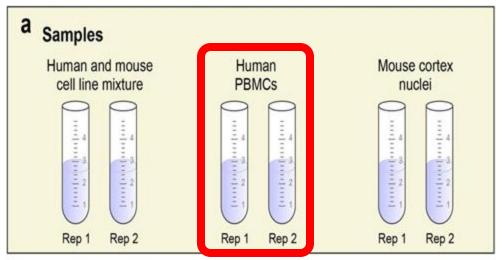
- Train on one dataset, evaluate on another
- More realistic scenario
- More challenging, data is not aligned

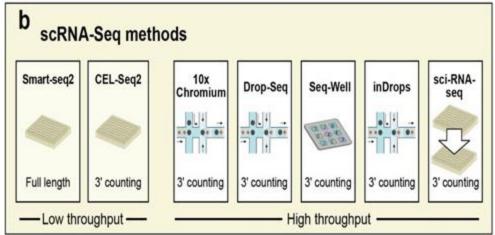






Experiment 2: inter-dataset evaluation

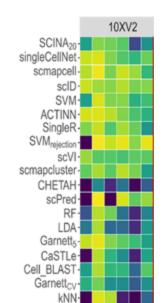








Prediction across protocols



LAmbDA-NMC-

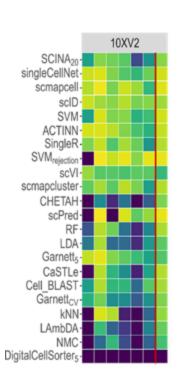
DigitalCellSorter₅-

Median F1-score

Training set 0 0.25 0.5 0.75



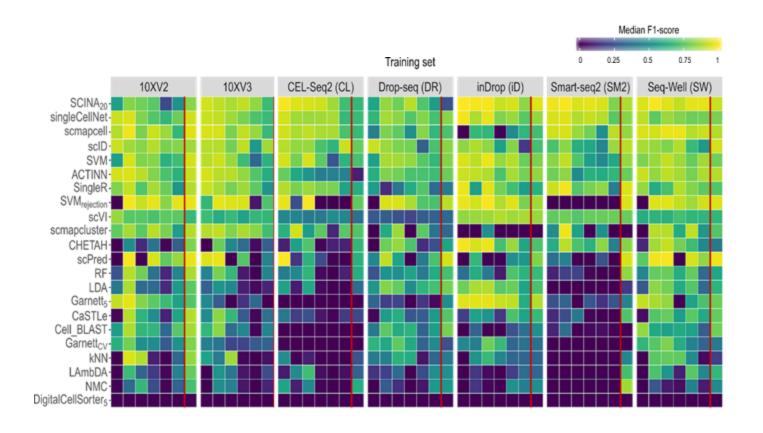
SM2 10XV3 CL DR iD SW



Training set

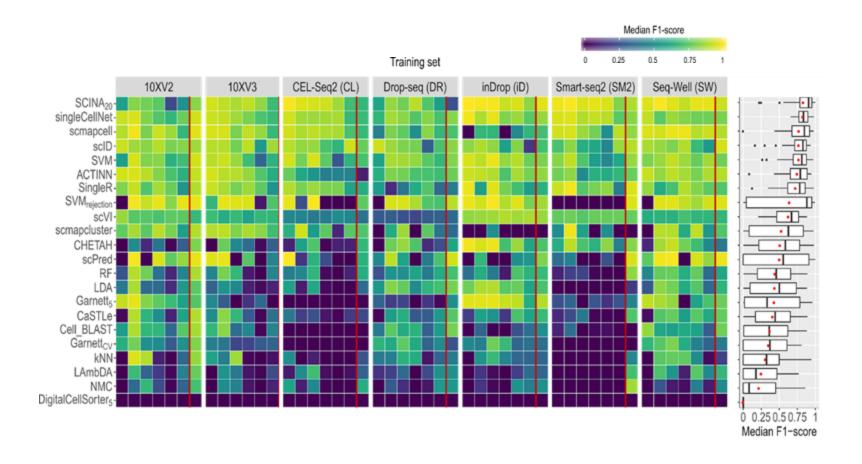








Test set SciLifeLab



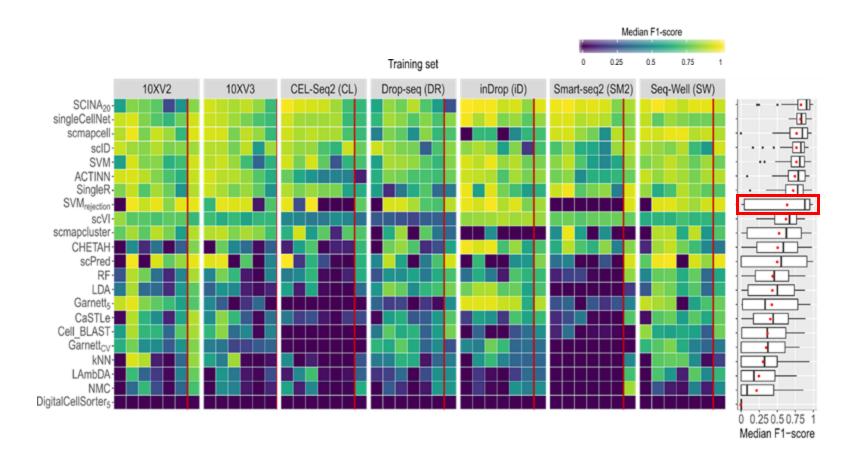


Test set

Test set

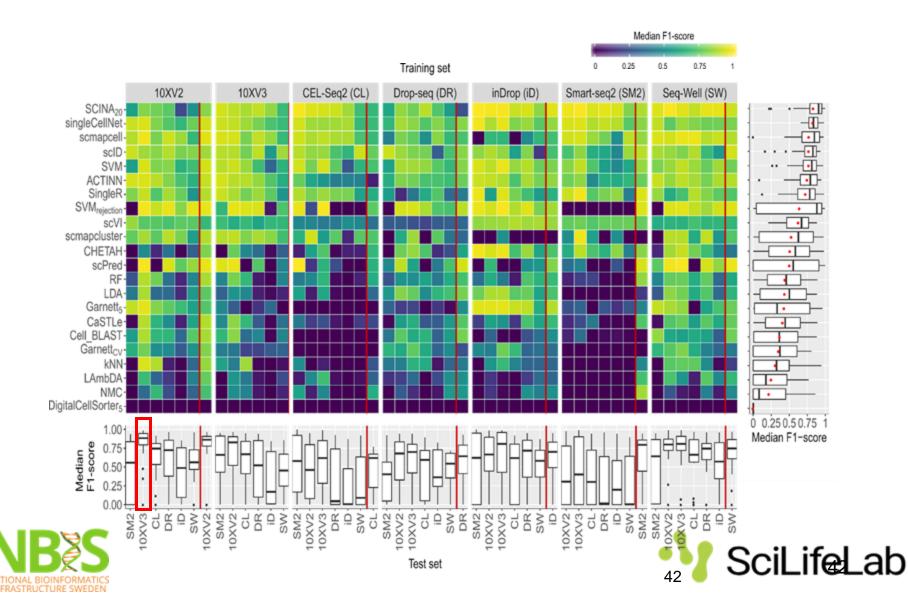
Test set

Test set





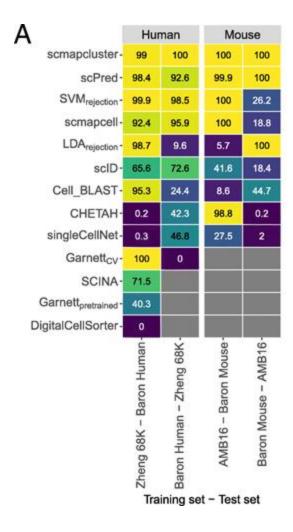


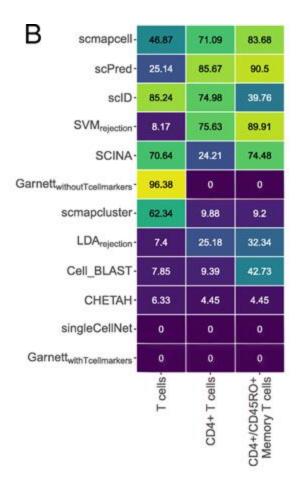


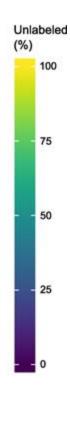




Experiment 3: rejection evaluation





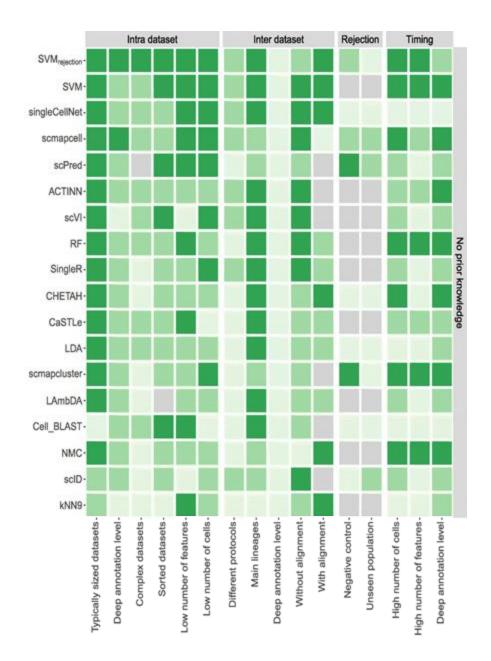


Cell population





Performance Summary



- Good
- Intermediate
- Poor





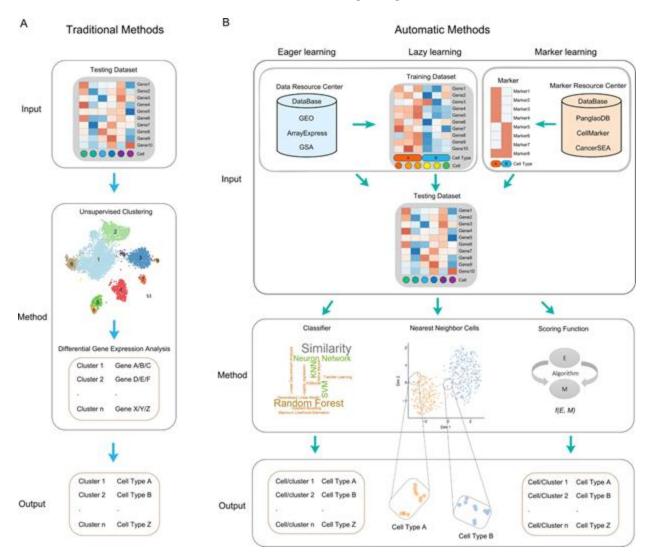
Conclusions so far

- 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: <u>https://github.com/tabdelaal/scRNAseg_Benchmark/</u>





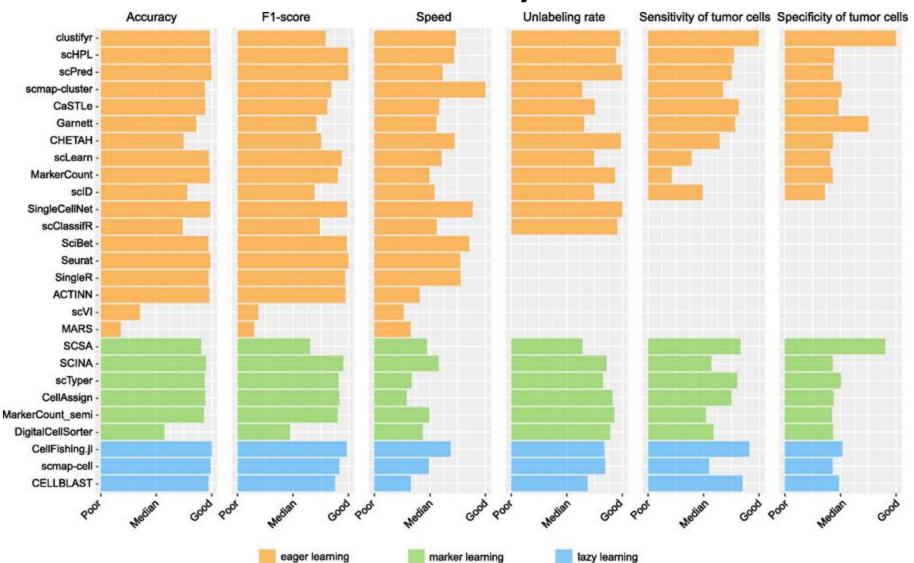
Benchmark paper 2021







Summary







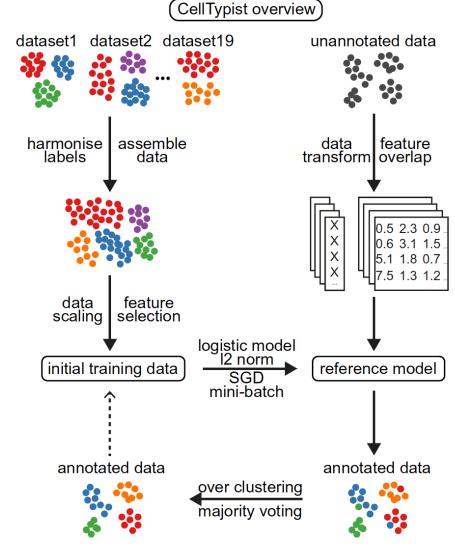
More recent methods/approaches?





CellTypist

- 20 different tissues from 19 datasets
- Immune cells across different organs
- Logistic regression with stochastic gradient descent learning

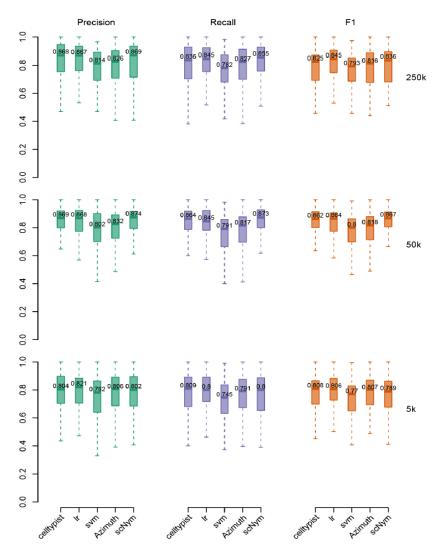






CellTypist

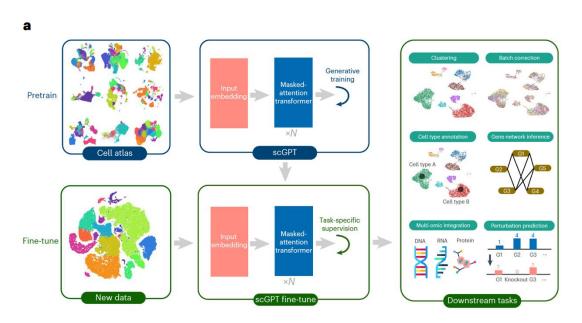
Benchmarking with other label-transferring methods

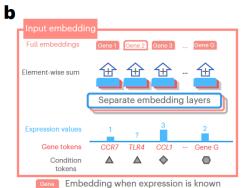




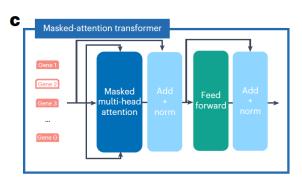


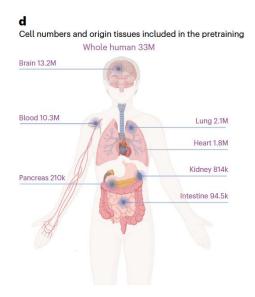
Generative learning is the next big thing? scGPT





Embedding when expression is unknown

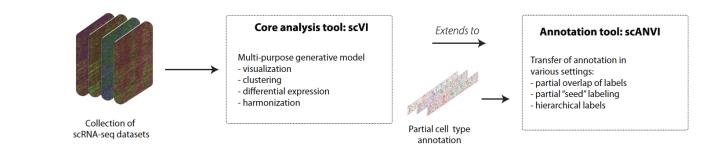


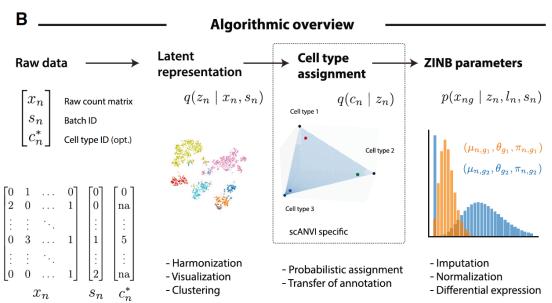


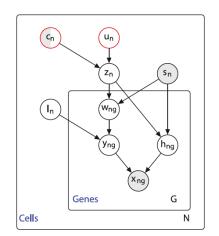




Generative learning is the next big thing? scANVI





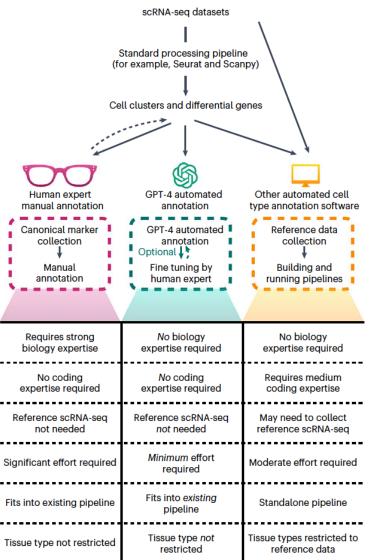


Probabilistic graphical model





GPTCelltype



nature methods



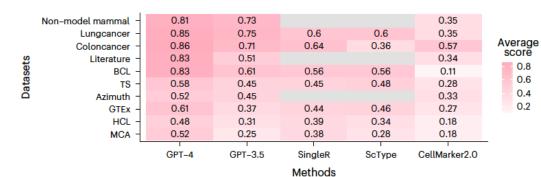
Brief Communication

https://doi.org/10.1038/s41592-024-02235-4

Assessing GPT-4 for cell type annotation in single-cell RNA-seq analysis

Received: 16 April 2023 Wenpin Hou ® ¹ ⊠ & Zhicheng Ji ® ² ⊠

Accepted: 5 March 2024







Challenges of automated cell type annotation?

- There is no ground truth for cell type annotation within a specific dataset
- Biology is complex, and cell states vary continuously
- Delineations between cell types are imprecise
- Cancer cell annotation is even more challenging due to their heterogeneity
- It is crucial that annotation methods highlight areas of uncertainty that may require manual scrutiny
- Any solution?





popular Vote (popV)

nature genetics

Article

https://doi.org/10.1038/s41588-024-01993-3

9

Consensus prediction of cell type labels in single-cell data with popV

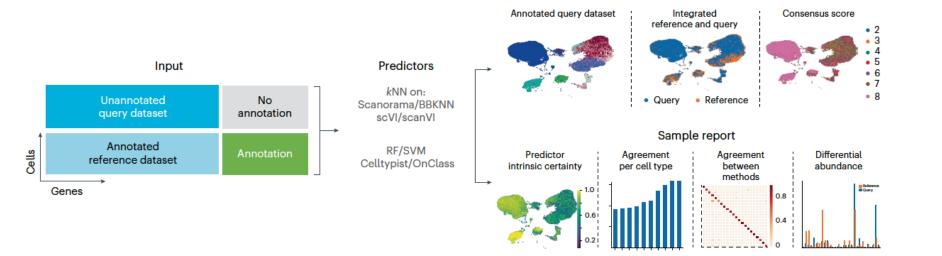
Received: 2 September 2023

Can Ergen ● ^{1,2}, Galen Xing ^{1,3}, Chenling Xu¹, Martin Kim², Michael Jayasuriya²,

Accepted: 18 October 2024

Published online: 20 November 2024

Result

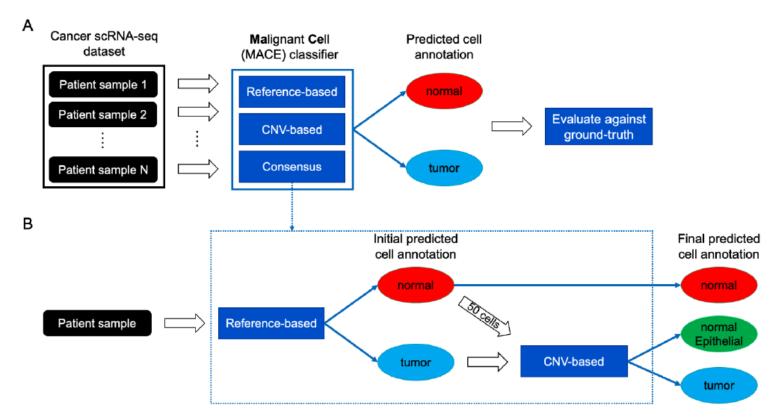






Malignant Cell (MACE) annotation

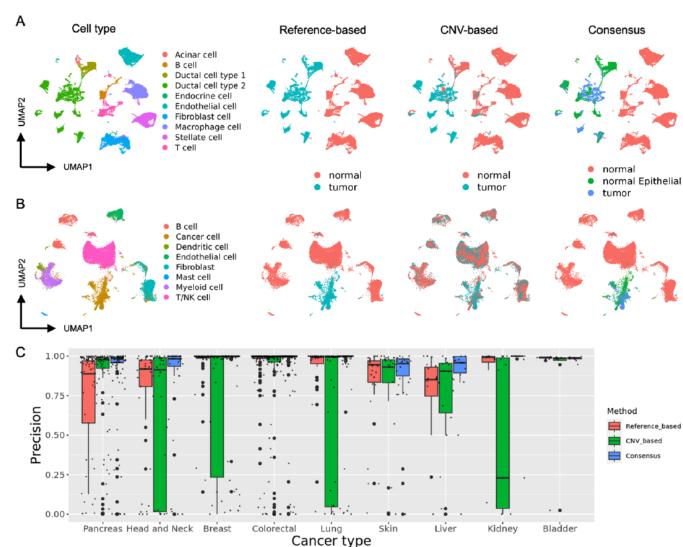
- CNV-based methods: SCEVAN
- Reference-based method: scATOMIC







Malignant Cell (MACE) annotation





Peng et al PDAC

Qian_et_al_Breast



Summary

- Cell identification is moving from unsupervised (clustering/visualization) to supervised (classification) learning
- Check what reference you are using!
 - The more similar the 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 make sense
- Consensus annotation methods combine the power of other automated methods and can be beneficial when addressing new, unknown cell types or disease samples.





Some useful resources

- Azimuth Seurat label transfer to reference sets
 - https://azimuth.hubmapconsortium.org/
 - online or R package
- DISCO CellMapper to several tissues (correlation-based)
 - <u>https://www.immunesinglecell.org/</u>
- Celltypist Regularised linear models with Stochastic Gradient Descent
 - <u>https://www.celltypist.org/</u>
 - online or python package
- scATOMIC (Random Forest) Pan-cancer TME cell type classifier
 - https://github.com/abelson-lab/scATOMIC





Acknowledgment

Most of the slides were adapted from previous presentations by Åsa Björklund (NBIS, scilifelab) Ahmed Mahfouz (Leiden University Medical Center, Björklund TU Delft)



