# *BTEP Presentation: scRNA – Cell Type Annotation*

Keyur Talsania
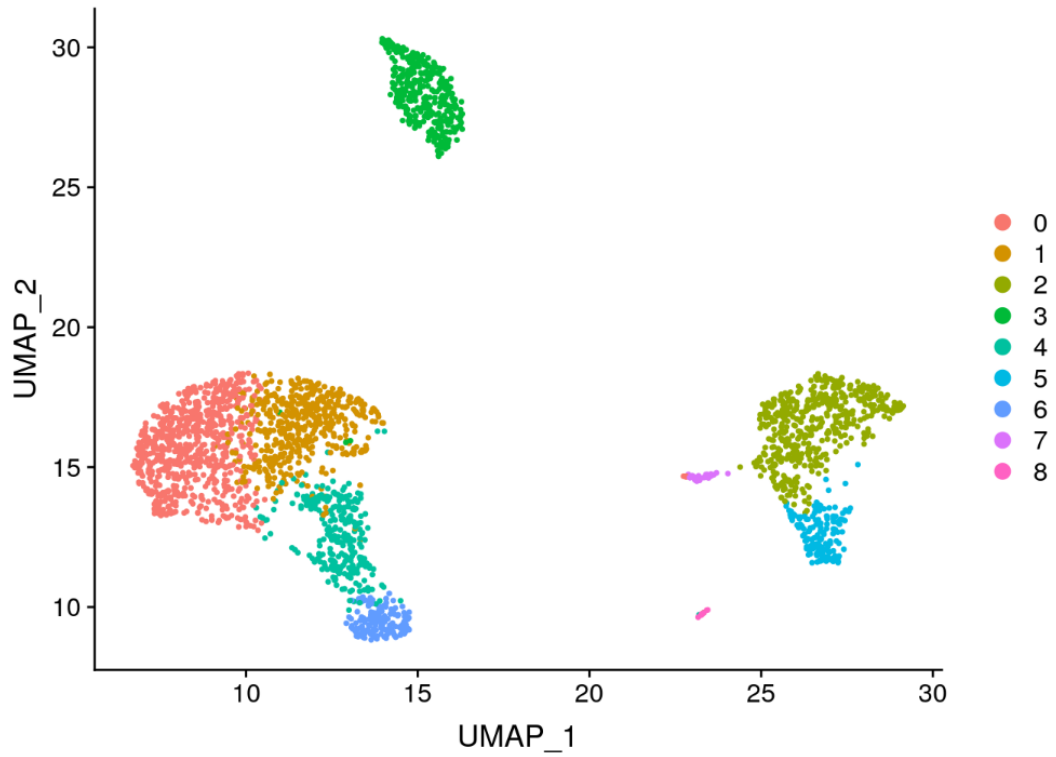
CCR-SF Bioinformatics Group
Advanced Biomedical and Computational Sciences

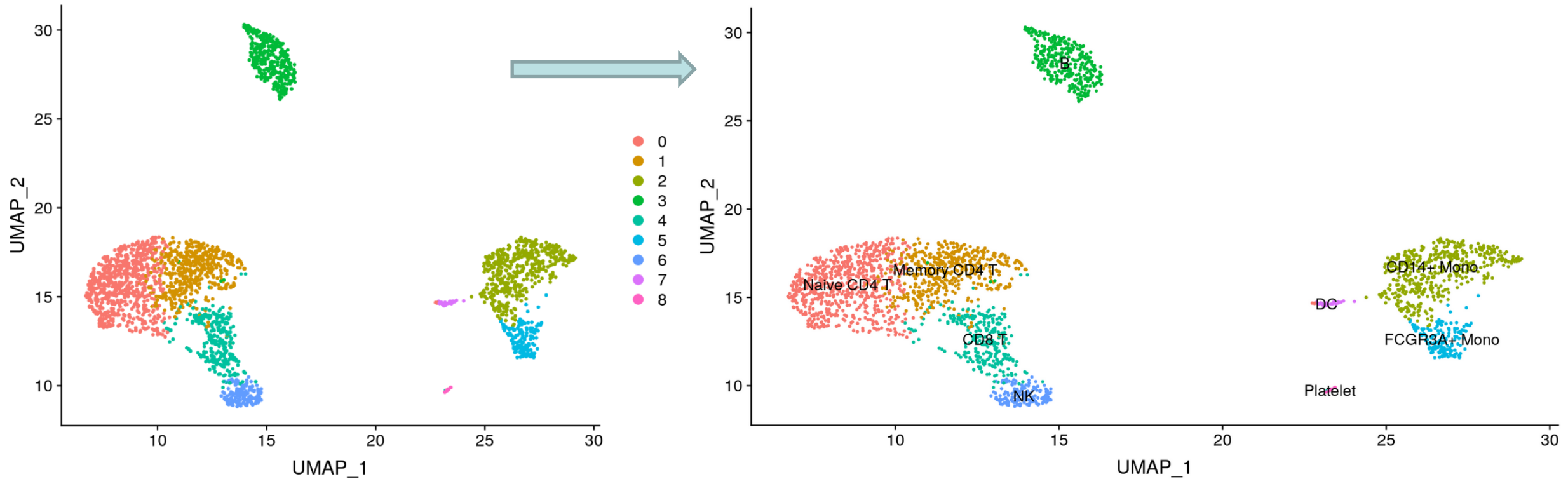Biomedical Informatics and Data Science (BIDS)  Directorate

Frederick National Laboratory for Cancer Research

# Outline

- **How it works**
  - Tools available
  - Publications

- **Tools we use**
  - SingleR, DigitalCellSorter
    - Label Transfer - SingleR

- **Concept of Automated Panels**
  - SF Pipeline's Automated Panels

# Why we need cell type annotation?

# Tools Available - https://www.scrna-tools.org

- ACTINN
- AltAnalyze
- bigSCale2
- cardelino
- CaSTLe
- celaref
- Cell-BLAST
- cellassign
- CellFishing
- CellO
- CHETAH
- ClusterMap

- DigitalCellSorter
- DistMap
- DropLasso
- FateID
- Garnett
- hscScore
- LAmbDA
- matchSCore2
- MetaNeighbor
- MIMOSCA
- Moana
- ParaDPMM

- scdney
- scID
- SCINA
- scmap
- scMatch
- scMCA
- scPred
- scVI
- SingleCellNet
- SingleR

# Types of Cell Annotation Tools

- **Supervised methods :**

  – which require a training dataset labeled with the corresponding cell populations in order to train the classifier

  – SingleR, SVM, ACTINN, scPred, CaSTle


- **Prior-knowledge based methods:**

  – for which either a marker gene file is required as an input or a pretrained classifier for specific cell populations is provided.

  – Garnett, DigitalCelllSorter, Moana

# Articles of comparison of Celltype Annotation tools

Check for updates

RESEARCH ARTICLE

REVISED **Evaluation of methods to assign cell type labels to cell clusters from single-cell RNA-sequencing data [version 2; peer review: 1 approved, 2 approved with reservations]**

J. Javier Diaz-Mejia [1-3], Elaine C. Meng[3], Alexander R. Pico [4], Sonya A. MacParland [5-7], Troy Ketela[1], Trevor J. Pugh[1,8,9], Gary D. Bader [2,10], John H. Morris [3]

Genome Biology

**RESEARCH**                                                                                                  **Open Access**

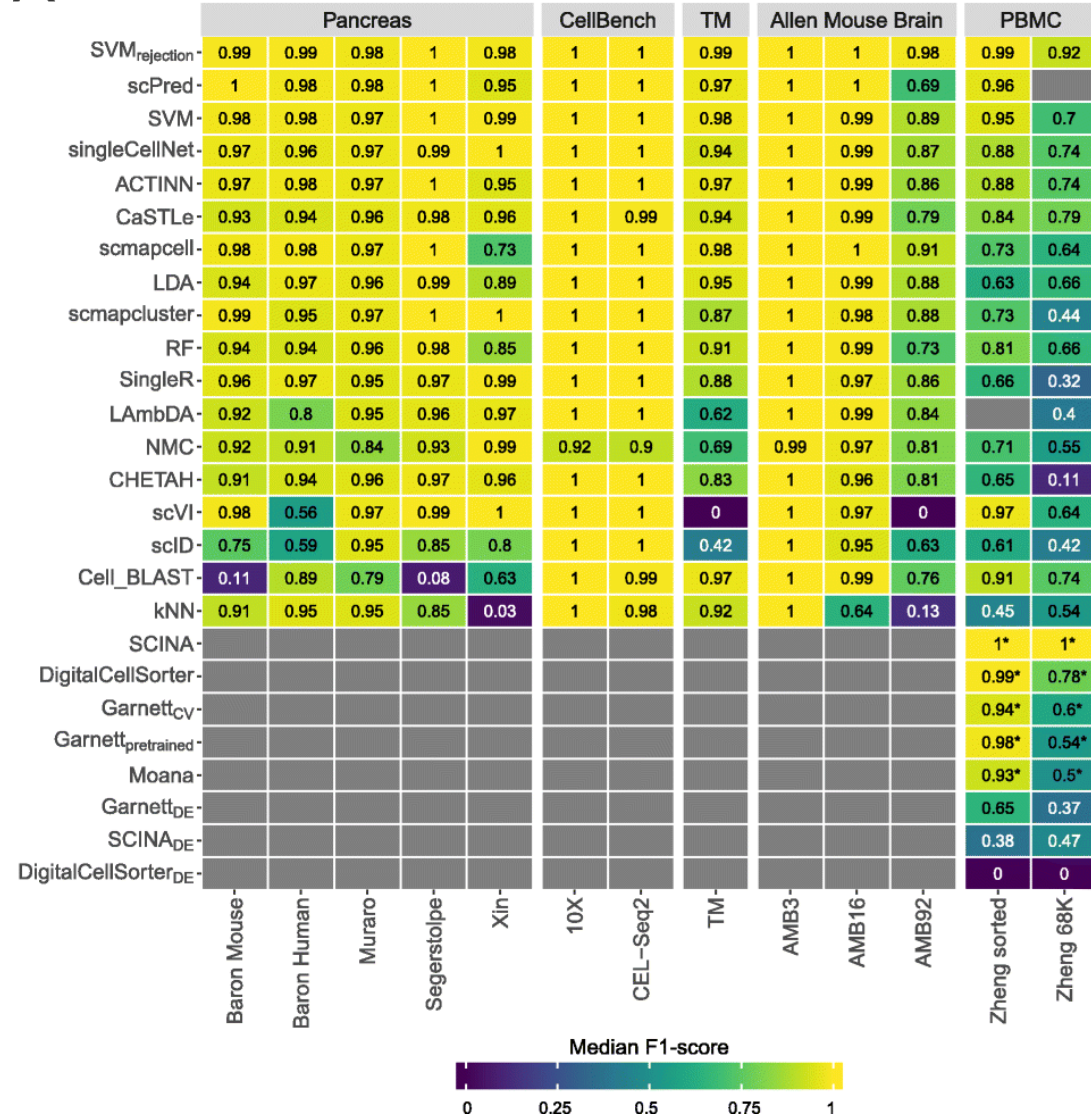# A comparison of automatic cell identification methods for single-cell RNA sequencing data

Check for updates

Tamim Abdelaal[1,2†], Lieke Michielsen[1,2†], Davy Cats[3], Dylan Hoogduin[3], Hailiang Mei[3], Marcel J. T. Reinders[1,2] and Ahmed Mahfouz[1,2*]

# Abdelaal et al.

| Name | Version | Language | Underlying classifier | Prior knowledge | Rejection option | Reference |
|------|---------|----------|----------------------|-----------------|------------------|-----------|
| Garnett | 0.1.4 | R | Generalized linear model | Yes | Yes | [14] |
| Moana | 0.1.1 | Python | SVM with linear kernel | Yes | No | [15] |
| DigitalCellSorter | GitHub version: e369a34 | Python | Voting based on cell type markers | Yes | No | [16] |
| SCINA | 1.1.0 | R | Bimodal distribution fitting for marker genes | Yes | No | [17] |
| scVI | 0.3.0 | Python | Neural network | No | No | [18] |
| Cell-BLAST | 0.1.2 | Python | Cell-to-cell similarity | No | Yes | [19] |
| ACTINN | GitHub version: 563bcc1 | Python | Neural network | No | No | [20] |
| LAmbDA | GitHub version: 3891d72 | Python | Random forest | No | No | [21] |
| scmapcluster | 1.5.1 | R | Nearest median classifier | No | Yes | [22] |
| scmapcell | 1.5.1 | R | kNN | No | Yes | [22] |
| scPred | 0.0.0.9000 | R | SVM with radial kernel | No | Yes | [23] |
| CHETAH | 0.99.5 | R | Correlation to training set | No | Yes | [24] |
| CaSTLe | GitHub version: 258b278 | R | Random forest | No | No | [25] |
| SingleR | 0.2.2 | R | Correlation to training set | No | No | [26] |
| scID | 0.0.0.9000 | R | LDA | No | Yes | [27] |
| singleCellNet | 0.1.0 | R | Random forest | No | No | [28] |
| LDA | 0.19.2 | Python | LDA | No | No | [29] |
| NMC | 0.19.2 | Python | NMC | No | No | [29] |
| RF | 0.19.2 | Python | RF (50 trees) | No | No | [29] |
| SVM | 0.19.2 | Python | SVM (linear kernel) | No | No | [29] |
| SVM$_{rejection}$ | 0.19.2 | Python | SVM (linear kernel) | No | Yes | [29] |
| kNN | 0.19.2 | Python | kNN ($k = 9$) | No | No | [29] |

# Abdelaal et al.

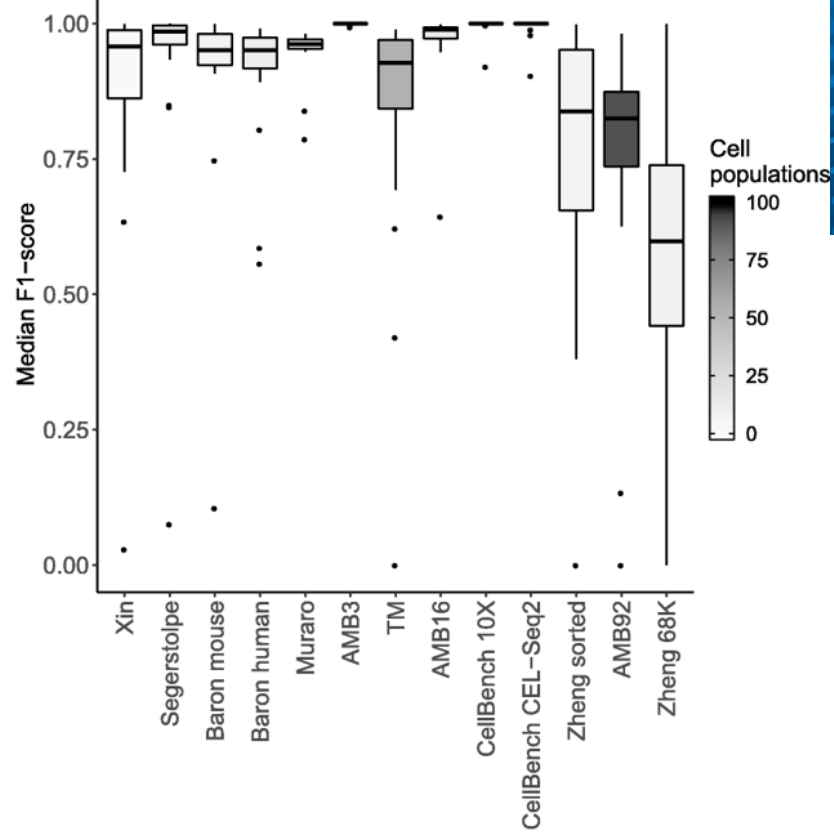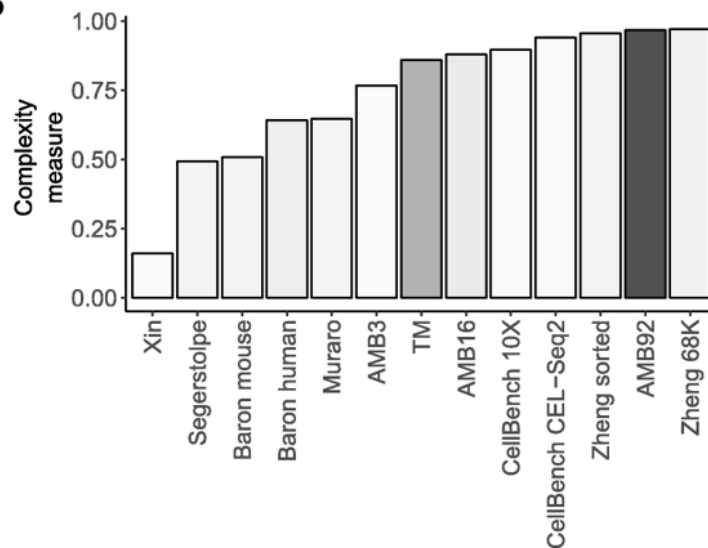| Dataset | No. of cells | No. of genes | No. of cell populations (>10 cells) | Description | Protocol | Reference |
|---|---|---|---|---|---|---|
| Baron (Mouse)[a] | 1886 | 14,861 | 13 (9) | Mouse pancreas | inDrop | [30] |
| Baron (Human)[a,b] | 8569 | 17,499 | 14 (13) | Human pancreas | inDrop | [30] |
| Muraro[a,b] | 2122 | 18,915 | 9 (8) | Human pancreas | CEL-Seq2 | [31] |
| Segerstolpe[a,b] | 2133 | 22,757 | 13 (9) | Human pancreas | SMART-Seq2 | [32] |
| Xin[a,b] | 1449 | 33,889 | 4 (4) | Human pancreas | SMARTer | [33] |
| CellBench 10X[a,b] | 3803 | 11,778 | 5 (5) | Mixture of five human lung cancer cell lines | 10X chromium | [34] |
| CellBench CEL-Seq2[a,b] | 570 | 12,627 | 5 (5) | Mixture of five human lung cancer cell lines | CEL-Seq2 | [34] |
| TM[a] | 54,865 | 19,791 | 55 (55) | Whole *Mus musculus* | SMART-Seq2 | [6] |
| AMB[a] | 12,832 | 42,625 | 4/22/110 (3/16/92) | Primary mouse visual cortex | SMART-Seq v4 | [35] |
| Zheng sorted[a] | 20,000 | 21,952 | 10 (10) | FACS-sorted PBMC | 10X CHROMIUM | [36] |
| Zheng 68K[a] | 65,943 | 20,387 | 11 (11) | PBMC | 10X CHROMIUM | [36] |
| VISp[b] (Mouse) | 12,832 | 42,625 | 3/36 (3/34) | Primary visual cortex | SMART-Seq v4 | [35] |
| ALM[b] (Mouse) | 8758 | 42,461 | 3/37 (3/34) | Anterior lateral motor area | SMART-Seq v4 | [35] |
| MTG[b] (Human) | 14,636 | 16,161 | 3/35 (3/34) | Middle temporal gyrus | SMART-Seq v4 | [37] |
| PbmcBench pbmc1.10Xv2[b] | 6444 | 33,694 | 9 (9) | PBMC | 10X version 2 | [38] |
| PbmcBench pbmc1.10Xv3[b] | 3222 | 33,694 | 8 (8) | PBMC | 10X version 3 | [38] |
| PbmcBench pbmc1.CL[b] | 253 | 33,694 | 7 (7) | PBMC | CEL-Seq2 | [38] |
| PbmcBench pbmc1.DR[b] | 3222 | 33,694 | 9 (9) | PBMC | Drop-Seq | [38] |
| PbmcBench pbmc1.iD[b] | 3222 | 33,694 | 7 (7) | PBMC | inDrop | [38] |
| PbmcBench pbmc1.SM2[b] | 253 | 33,694 | 6 (6) | PBMC | SMART-Seq2 | [38] |
| PbmcBench pbmc1.SW[b] | 3176 | 33,694 | 7 (7) | PBMC | Seq-Well | [38] |
| PbmcBench pbmc2.10Xv2[b] | 3362 | 33,694 | 9 (9) | PBMC | 10X version 2 | [38] |
| PbmcBench pbmc2.CL[b] | 273 | 33,694 | 5 (5) | PBMC | CEL-Seq2 | [38] |
| PbmcBench pbmc2.DR[b] | 3362 | 33,694 | 6 (6) | PBMC | Drop-Seq | [38] |
| PbmcBench pbmc2.iD[b] | 3362 | 33,694 | 9 (9) | PBMC | inDrop | [38] |
| PbmcBench pbmc2.SM2[b] | 273 | 33,694 | 6 (6) | PBMC | SMART-Seq2 | [38] |
| PbmcBench pbmc2.SW[b] | 551 | 33,694 | 4 (4) | PBMC | Seq-Well | [38] |

# Abdelaal et al. – All data

# Abdelaal et al. - Brain

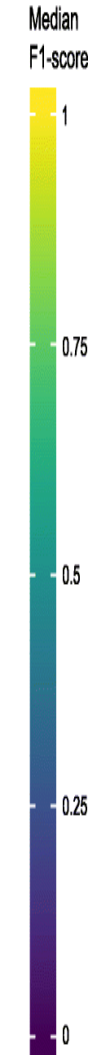# Abdelaal et al. - Pancreas

# Abdelaal et al.
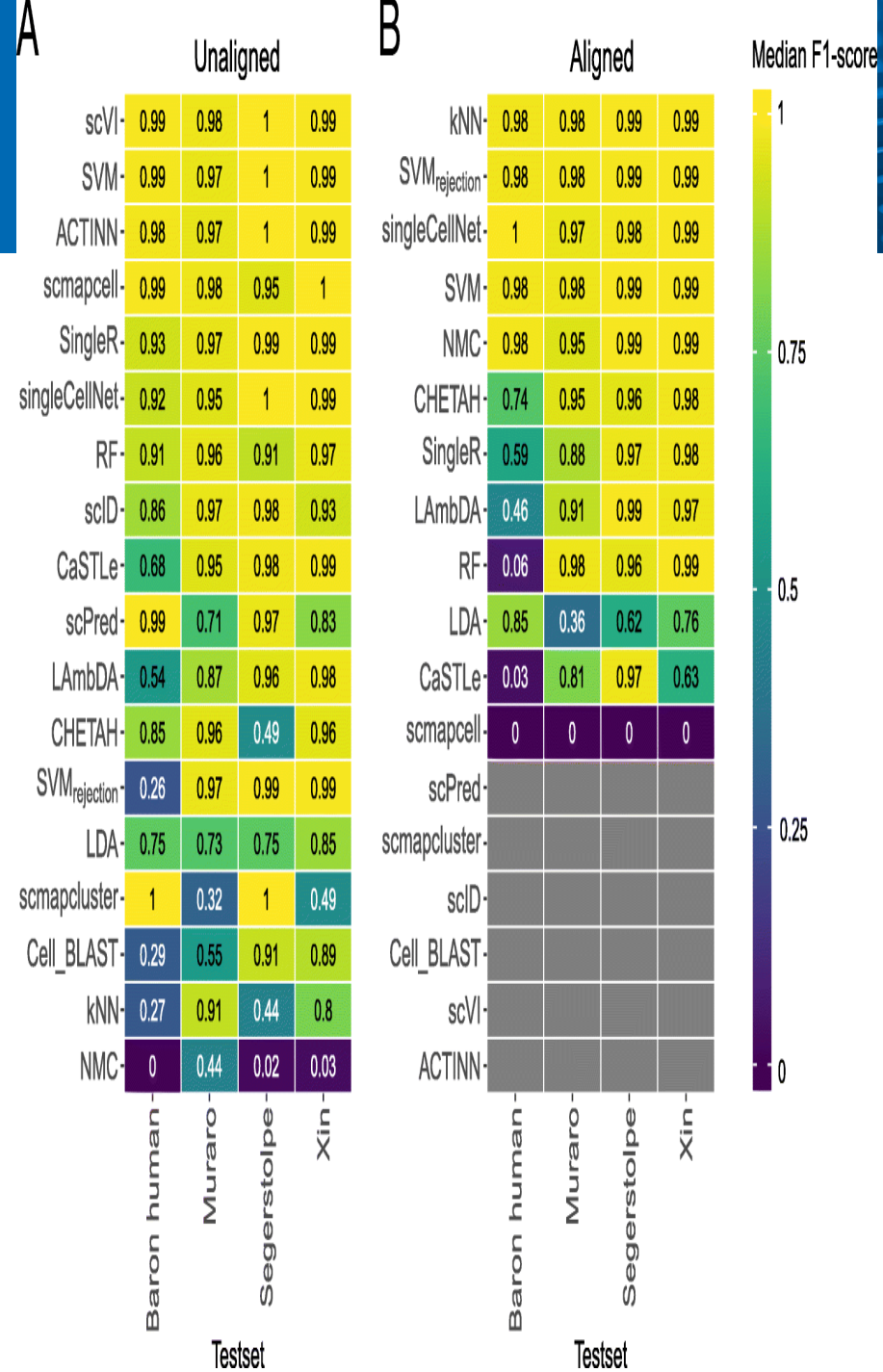
Frederick
National
Laboratory
for Cancer Research

Check for updates

**RESEARCH ARTICLE**

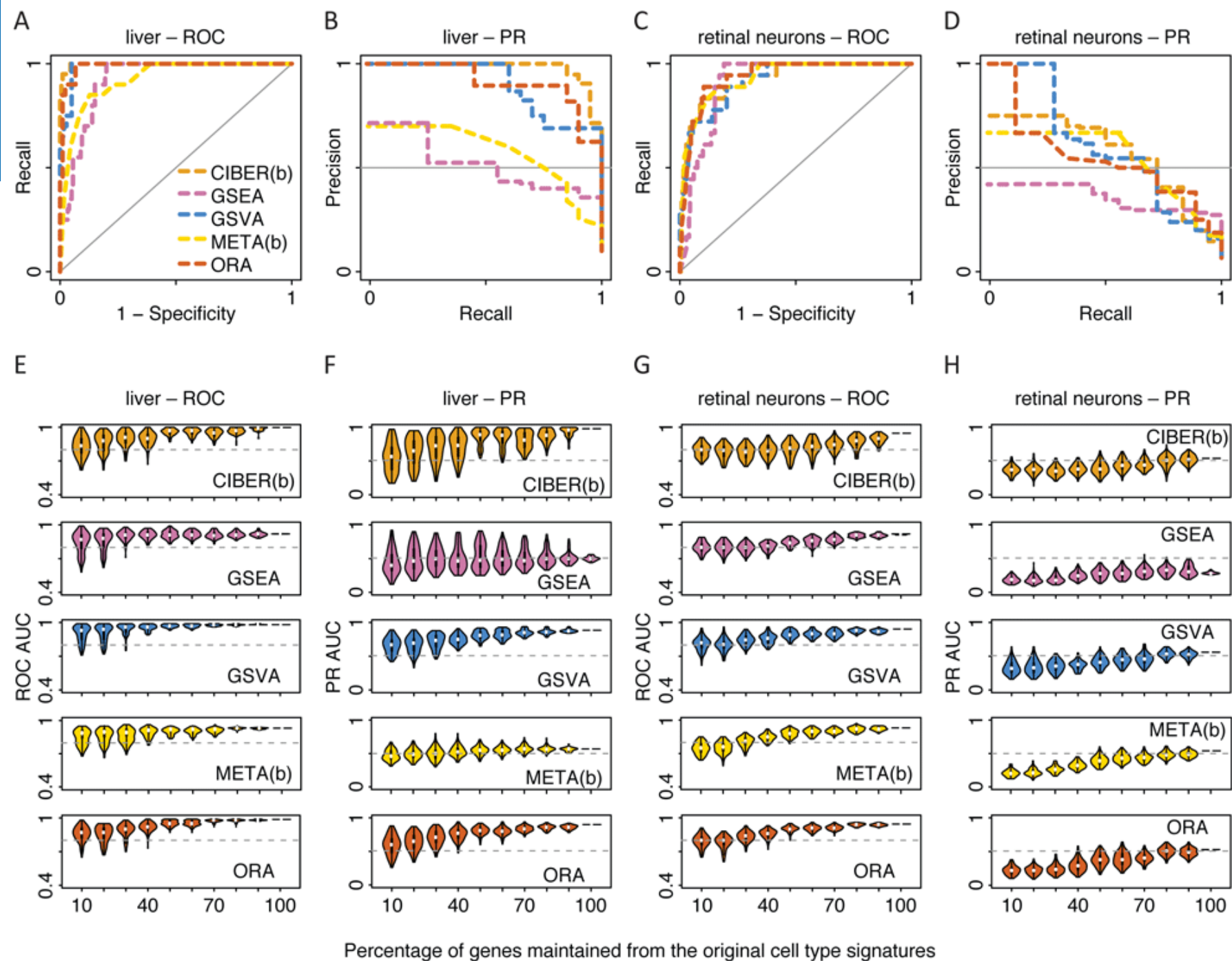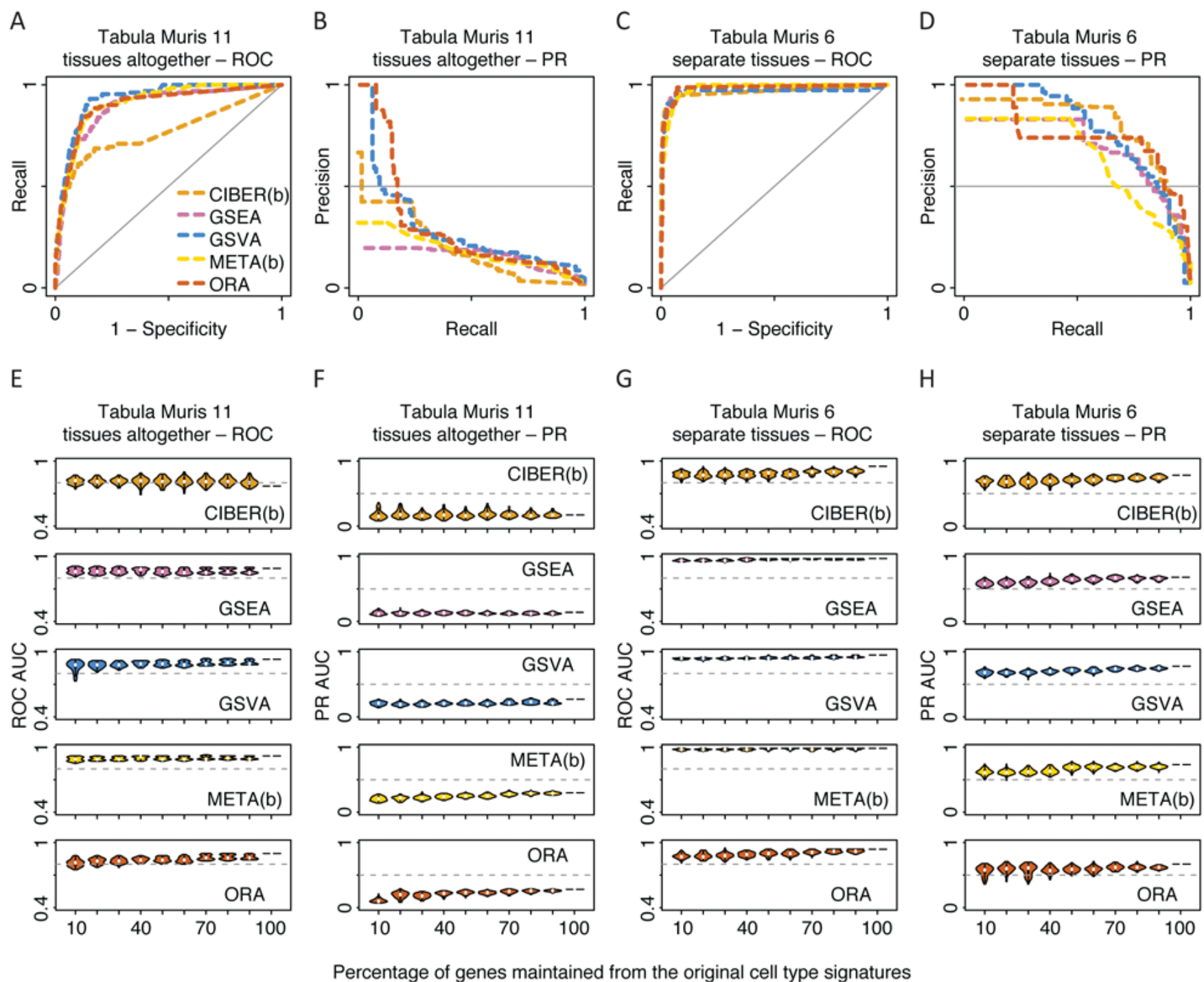**REVISED** **Evaluation of methods to assign cell type labels to cell clusters from single-cell RNA-sequencing data [version 2; peer review: 1 approved, 2 approved with reservations]**

J. Javier Diaz-Mejia [1-3], Elaine C. Meng[3], Alexander R. Pico [4], Sonya A. MacParland [5-7], Troy Ketela[1], Trevor J. Pugh[1,8,9], Gary D. Bader [2,10], John H. Morris [3]
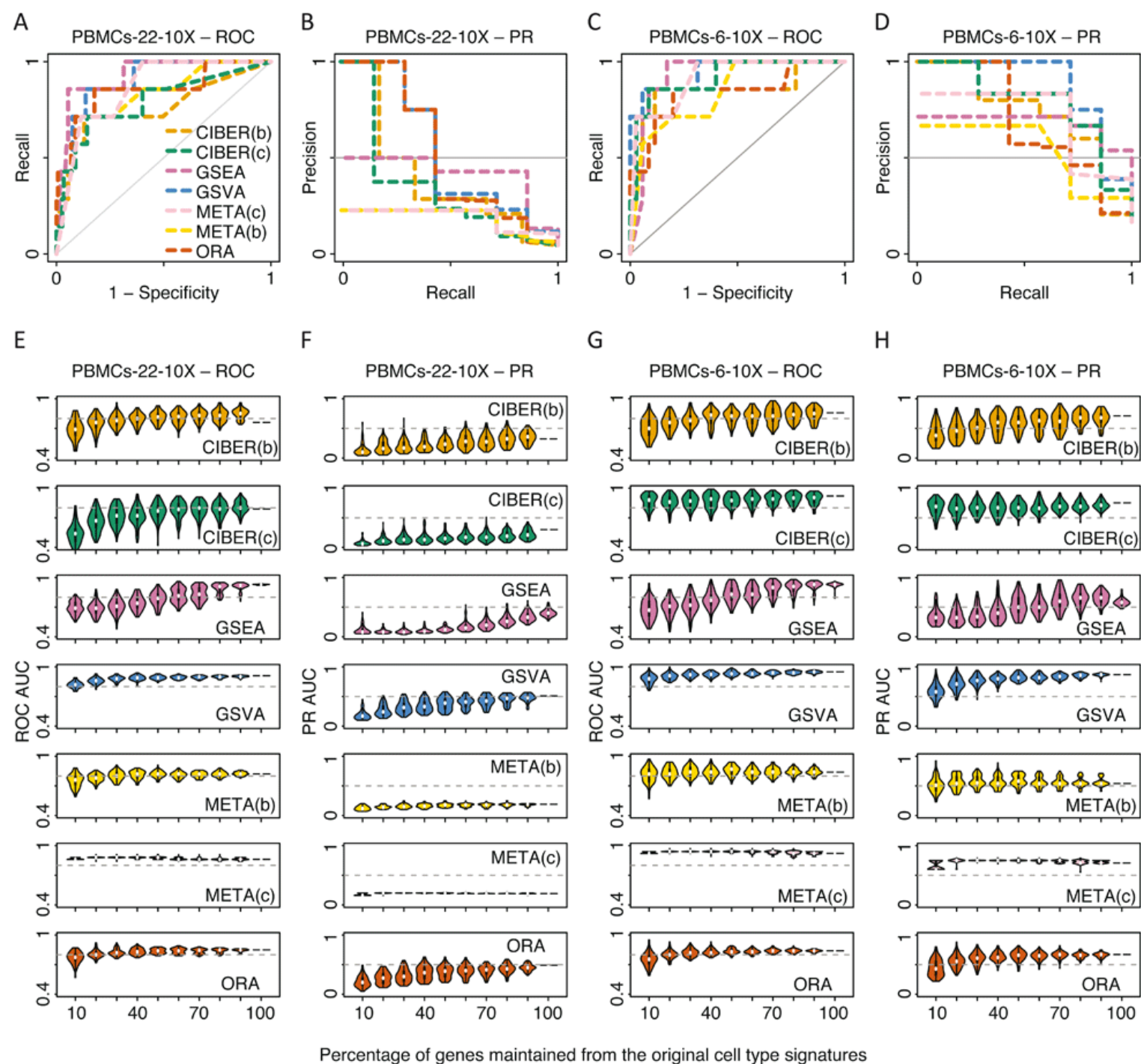
# Diaz-Mejia et al.

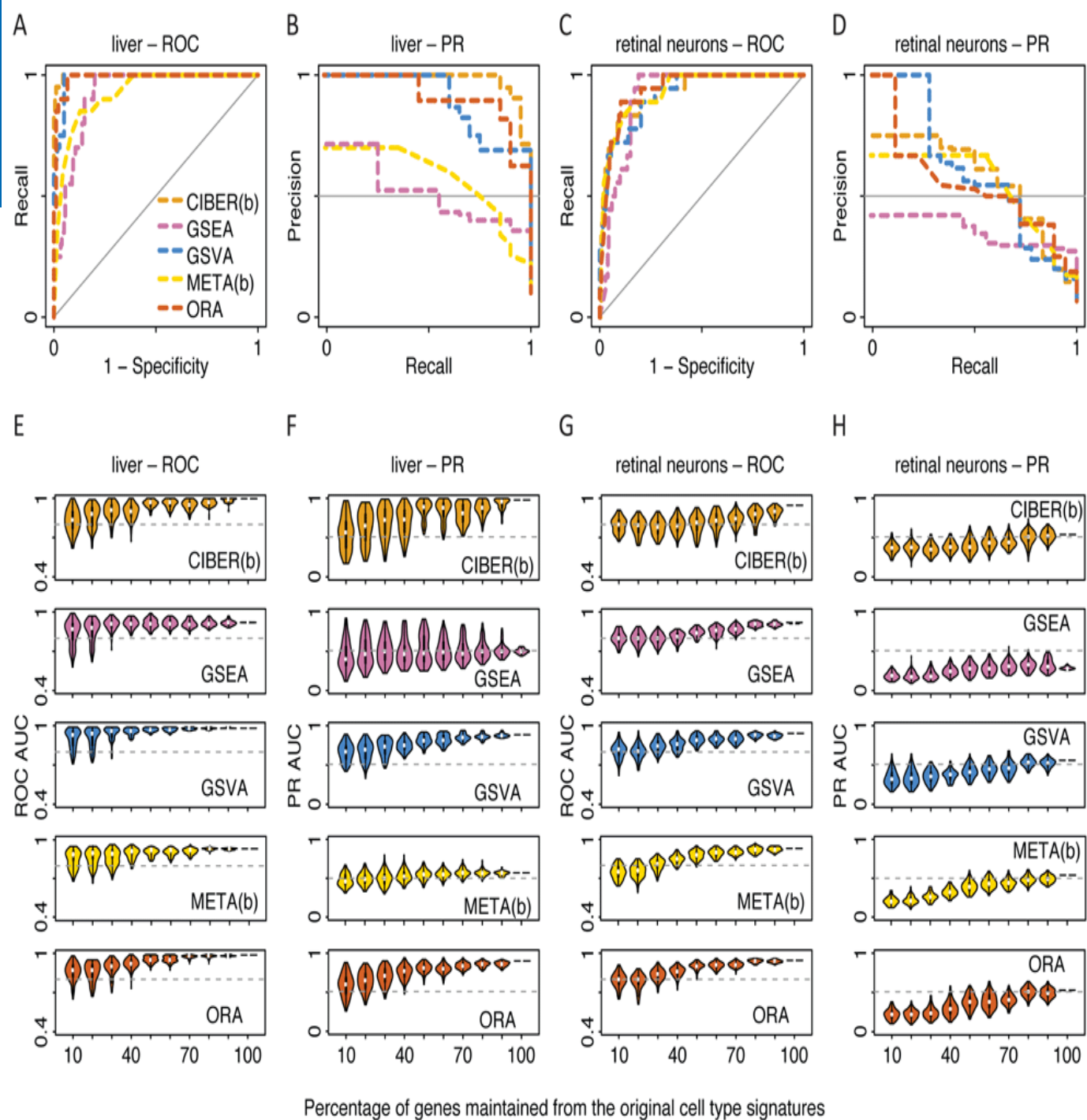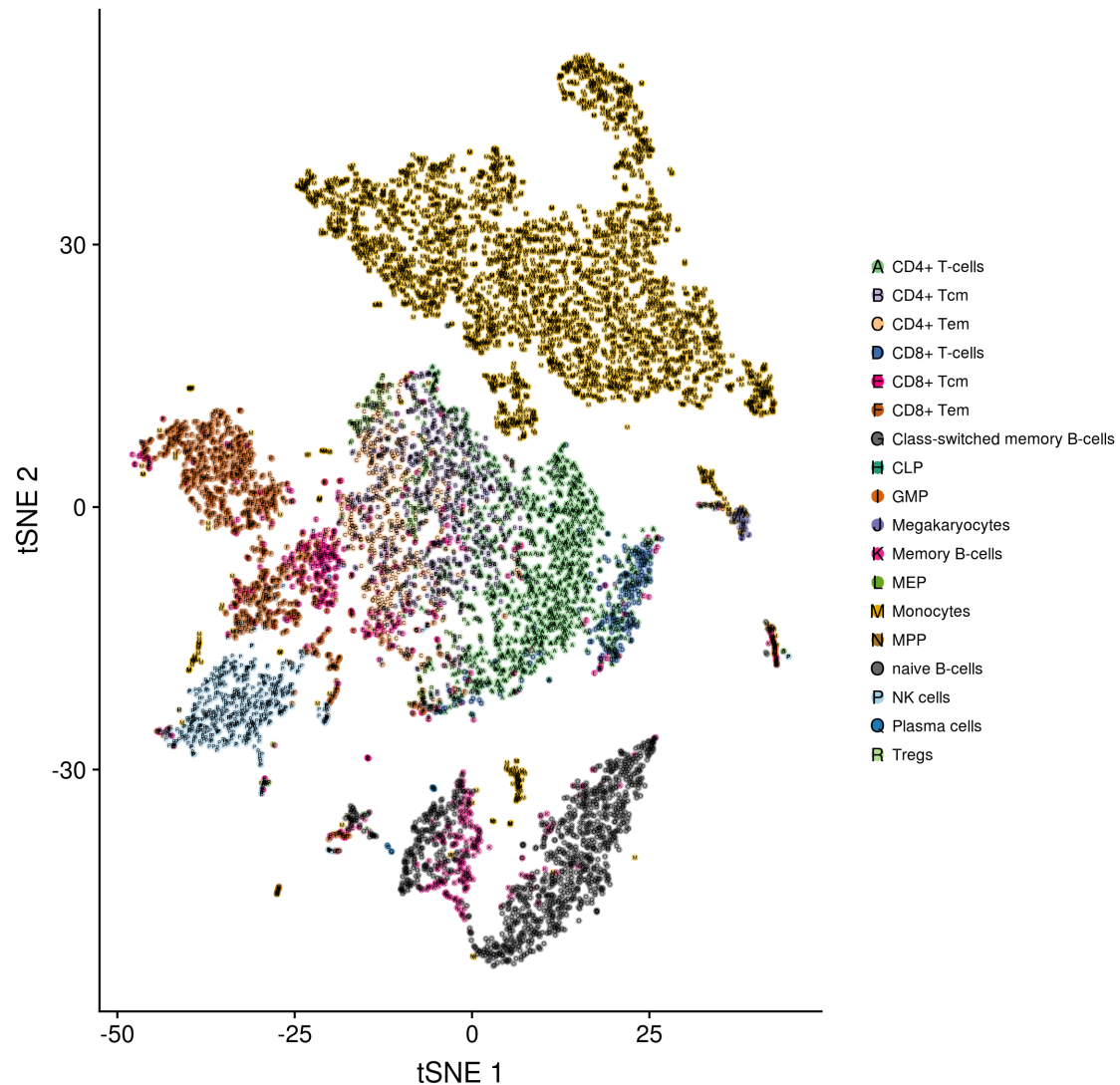| Acronym | Version | Name | Language | Reference |
|---------|---------|------|----------|-----------|
| CIBERSORT | 1.01 | Cell type Identification by Estimating Relative Subsets of RNA Transcripts | R and Java | (Newman *et al.*, 2015b) |
| GSEA | 3.0 | Gene Set Enrichment Analysis | Java | (Subramanian *et al.*, 2005) |
| GSVA | 1.30 | Gene Set Variation Analysis | R | (Hänzelmann *et al.*, 2013) |
| METANEIGHBOR | 1.3.1 | Meta-analysis via neighbor voting | R | (Crow *et al.*, 2018) |
| ORA | R( 3.5.1) | Over- representation Analysis | R | (Fisher, 1935; Goeman & Bühlmann, 2007) |

**Diaz-Mejia et al.**

Diaz-Mejia et al.

**Diaz-Mejia et al.**

# Diaz-Mejia et al.



Percentage of genes maintained from the original cell type signatures

# Tools we use….

# SingleR

| Data retrieval | Organism | Samples | Sample types | No. of main labels | No. of fine labels | Cell type focus |
|---|---|---|---|---|---|---|
| HumanPrimaryCellAtlasData() | human | 713 | microarrays of sorted cell populations | 37 | 157 | Non-specific |
| BlueprintEncodeData() | human | 259 | RNA-seq | 24 | 43 | Non-specific |
| DatabaseImmuneCellExpressionData() | human | 1561 | RNA-seq | 5 | 15 | Immune |
| NovershternHematopoieticData() | human | 211 | microarrays of sorted cell populations | 17 | 38 | Hematopoietic & Immune |
| MonacoImmuneData() | human | 114 | RNA-seq | 11 | 29 | Immune |
| ImmGenData() | mouse | 830 | microarrays of sorted cell populations | 20 | 253 | Hematopoietic & Immune |
| MouseRNAseqData() | mouse | 358 | RNA-seq | 18 | 28 | Non-specific |

Aran, Looney, Liu et al. Reference-based analysis of lung single-cell sequencing reveals a transitional profibrotic macrophage. Nature Immunology (2019)

# SingleR

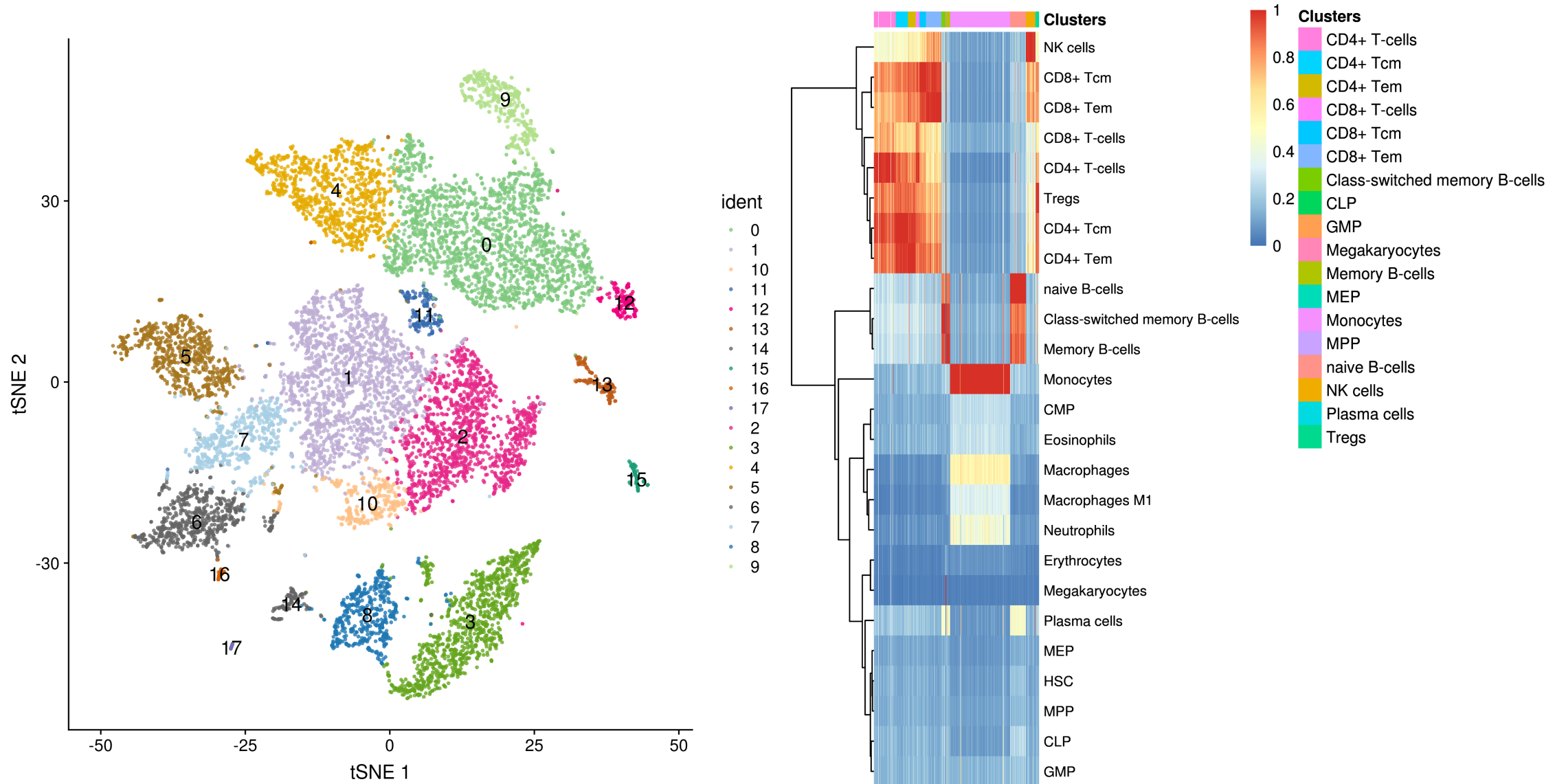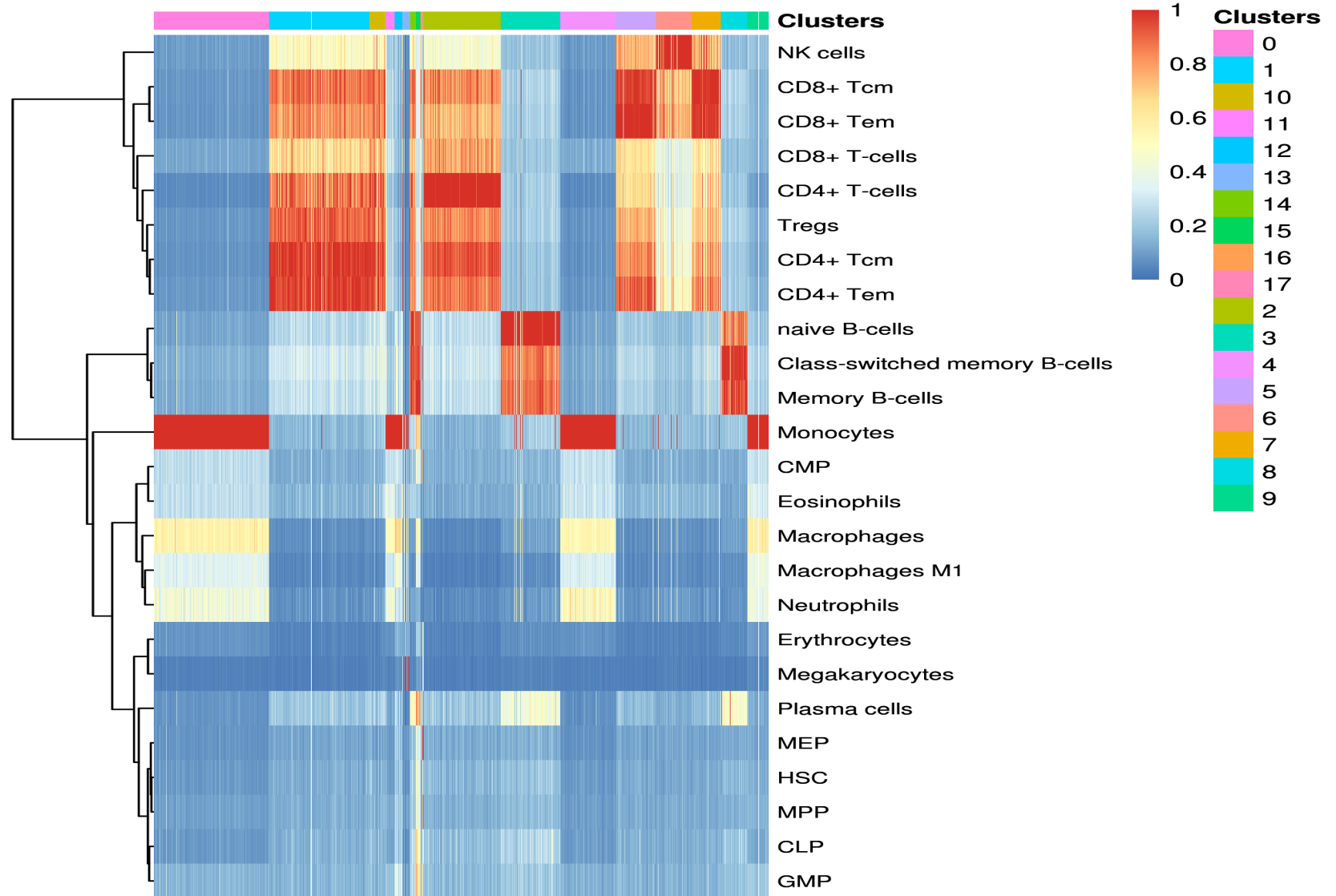| Data retrieval | Organism | Samples | Sample types | No. of main labels | No. of fine labels | Cell type focus |
|---|---|---|---|---|---|---|
| HumanPrimaryCellAtlasData() | human | 713 | microarrays of sorted cell populations | 37 | 157 | Non-specific |
| BlueprintEncodeData() | human | 259 | RNA-seq | 24 | 43 | Non-specific |
| DatabaseImmuneCellExpressionData() | human | 1561 | RNA-seq | 5 | 15 | Immune |
| NovershternHematopoieticData() | human | 211 | microarrays of sorted cell populations | 17 | 38 | Hematopoietic & Immune |
| MonacoImmuneData() | human | 114 | RNA-seq | 11 | 29 | Immune |
| ImmGenData() | mouse | 830 | microarrays of sorted cell populations | 20 | 253 | Hematopoietic & Immune |
| MouseRNAseqData() | mouse | 358 | RNA-seq | 18 | 28 | Non-specific |

+ Any dataset which is annotated!

# SingleR - PBMC



| Data retrieval | Organism | Samples | Sample types | No. of main labels | No. of fine labels | Cell type focus |
|---|---|---|---|---|---|---|
| HumanPrimaryCellAtlasData() | human | 713 | microarrays of sorted cell populations | 37 | 157 | Non-specific |
| BlueprintEncodeData() | human | 259 | RNA-seq | 24 | 43 | Non-specific |
| DatabaseImmuneCellExpressionData() | human | 1561 | RNA-seq | 5 | 15 | Immune |
| NovershternHematopoieticData() | human | 211 | microarrays of sorted cell populations | 17 | 38 | Hematopoietic & Immune |
| MonacoImmuneData() | human | 114 | RNA-seq | 11 | 29 | Immune |
| ImmGenData() | mouse | 830 | microarrays of sorted cell populations | 20 | 253 | Hematopoietic & Immune |
| MouseRNAseqData() | mouse | 358 | RNA-seq | 18 | 28 | Non-specific |

Legend:
- A CD4+ T-cells
- B CD4+ Tcm
- C CD4+ Tem
- D CD8+ T-cells
- E CD8+ Tcm
- F CD8+ Tem
- G Class-switched memory B-cells
- H CLP
- I GMP
- J Megakaryocytes
- K Memory B-cells
- L MEP
- M Monocytes
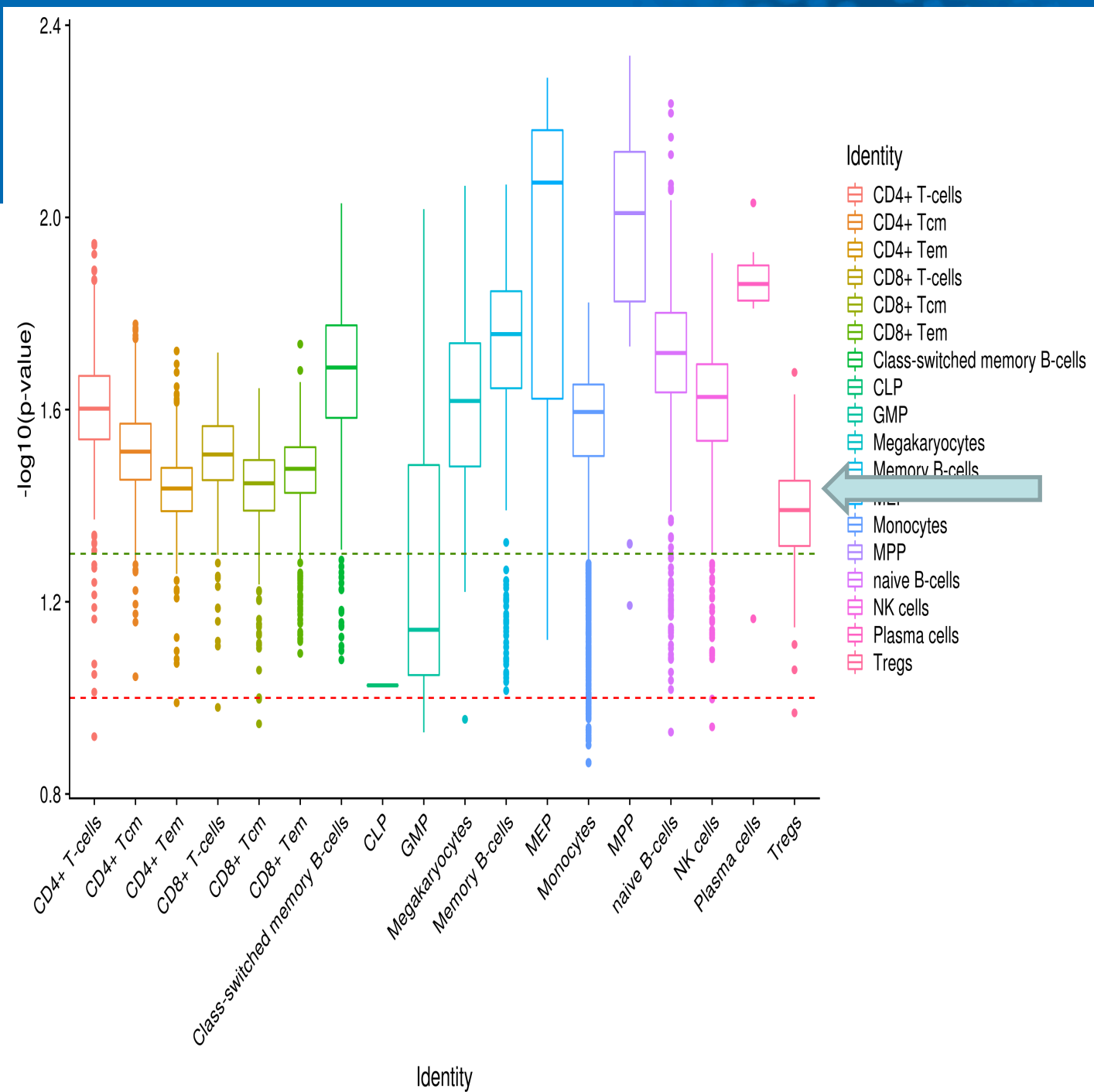- N MPP
- O naive B-cells
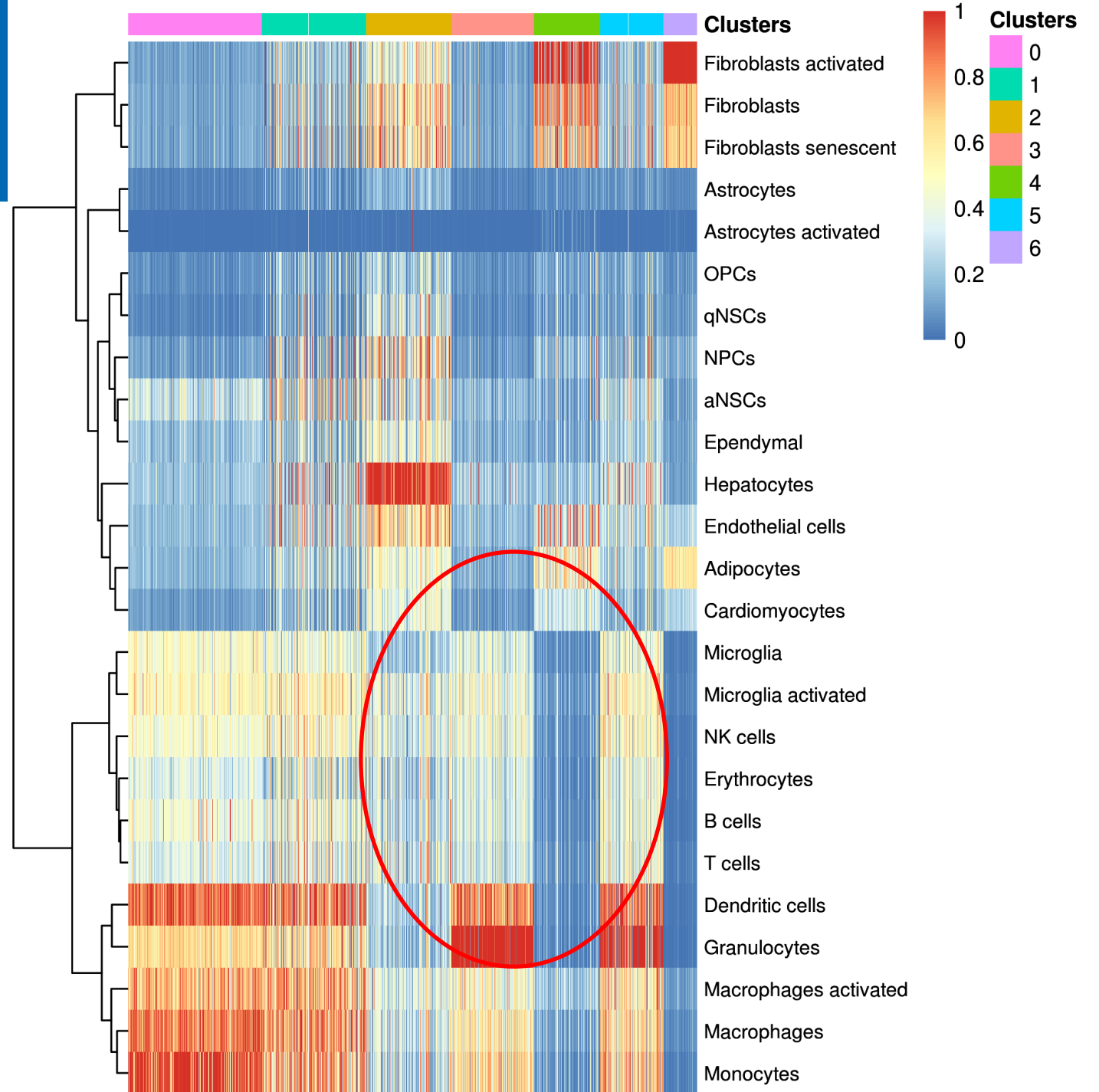- P NK cells
- Q Plasma cells
- R Tregs

# SingleR - PBMC

# SingleR - PBMC

# SingleR - Pancreas

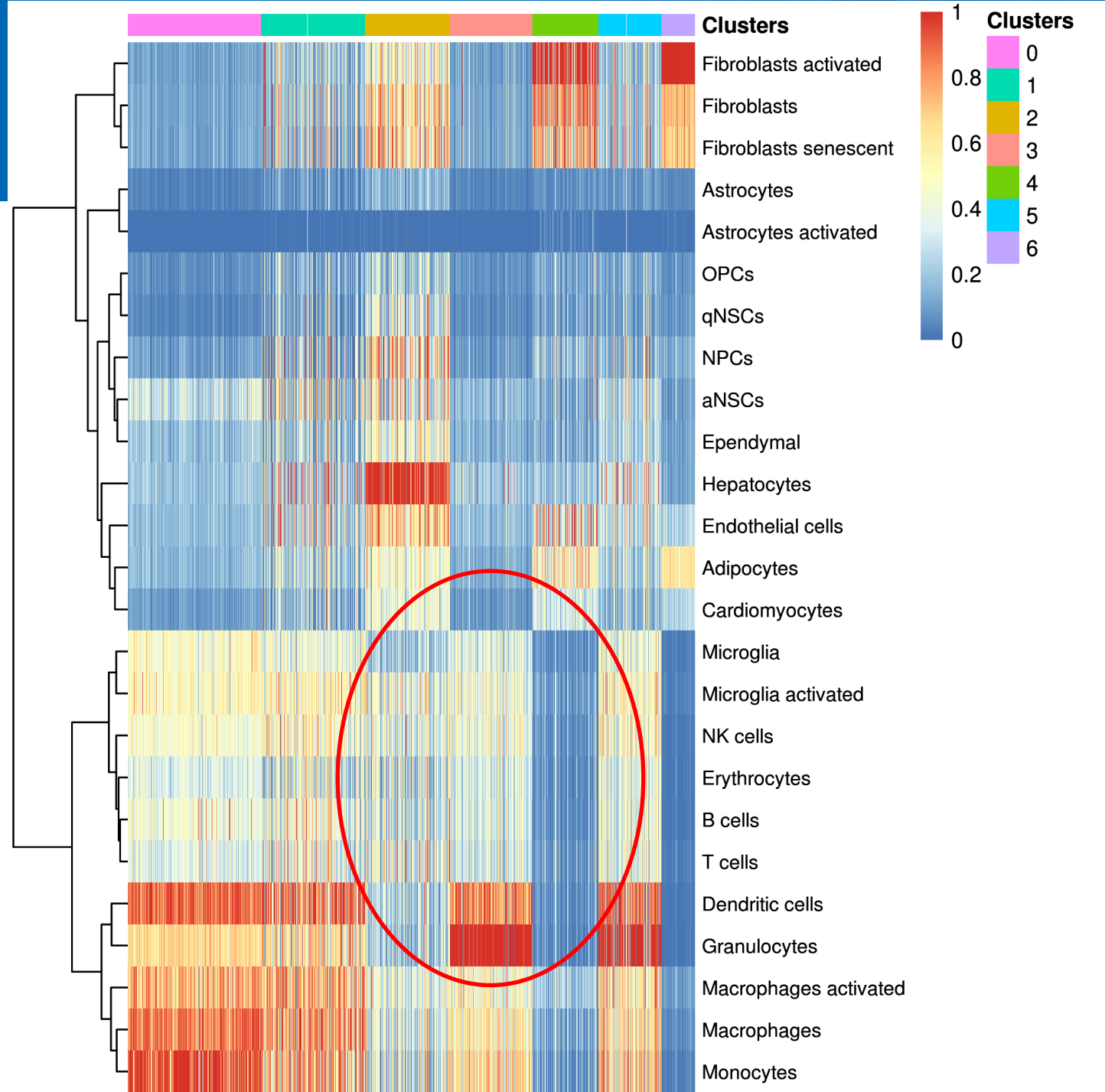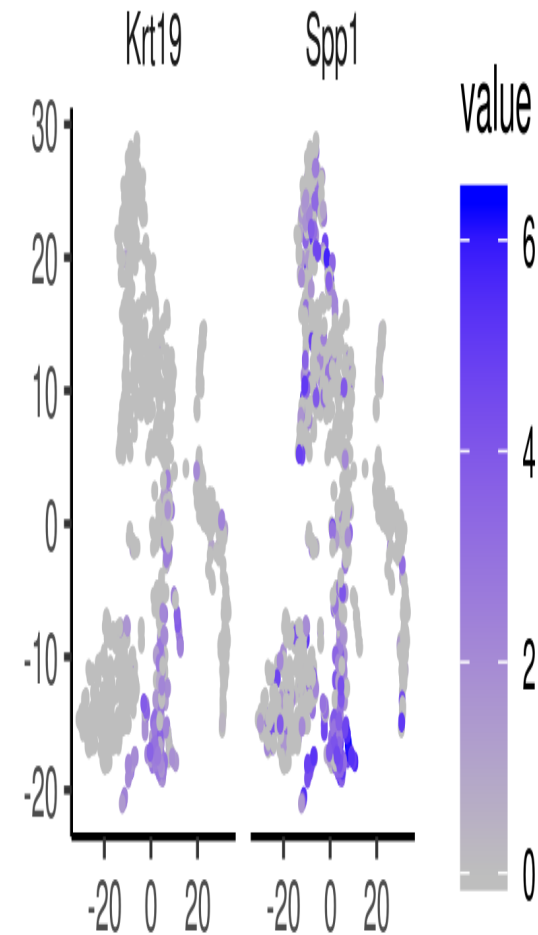SingleR - Pancreas

SingleR - Pancreas

# DigitalCellSorter

Sergii Domanskyi, Anthony Szedlak,
Nathaniel T Hawkins, Jiayin Wang,
Giovanni Paternostro & Carlo
Piermarocchi, *BMC Bioinformatics* volume
20, Article number: 369 (**2019**)

**1. Pre-preprocessing** of single cell mRNA sequencing data (gene expression data)
1. Cleaning: filling in missing values, zemoving all-zero genes and cells, converting gene index to a desired convention, etc.
2. Normalizing: rescaling all cells expression, log-transforming, etc.

**2. Quality control**

**3. Batch effects correction**

**4. Cells anomaly score evaluation**

**5. Dimensionality reduction**

**6. Clustering** (Hierarchical, K-Means, knn-graph-based, etc.)

**7. Annotating cell types**

**8. Vizualization**
1. t-SNE layout plot
2. Quality Control histogram plot
3. Marker expression t-SNE subplot
4. Marker-centroids expression plot
5. Voting results matrix plot
6. Cell types stacked barplot
7. Anomaly scores plot
8. Histogram null distribution plot
9. New markers plot
10. Sankey diagram (a.k.a. river plot)

**9. Post-processing** functions, e.g. extract cells of interest, find significantly expressed genes, plot marker expression of the cells of interest, etc.

## Pre-processing

Pre-filtering: keep only genes $i$ where $\sum_j X_{ij} > 0$

Normalization: $X_{ij}/\sum_{i'} X_{i'j}$

Transformation: $X_{ij} = \begin{cases} \log_2 X_{ij}, & X_{ij} > 0 \\ \log_2 m, & X_{ij} \leq 0 \end{cases}$

Post-filtering: keep only genes $i$ for which $\sigma_i/\langle\sigma\rangle$

## Clustering

Dimensionality reduction: PCA on $X_{ij}$ to 100 feat

Number of clusters: ARI on the PCA-reduced dat

Clustering PCA-reduced data: e.g. by Agglomerat clustering

## Cell type assignment

Prepare marker/cell type matrix $M_{km}$

Generate $P_{kc}(V_{kc})$ distribution

Calculate $\Lambda_{kc}$: z-score of $V_{kc}$ in

Get $T_c = \operatorname{argmax} \Lambda_{kc}$

## Plots for visua

Project PCA-reduced data on 2

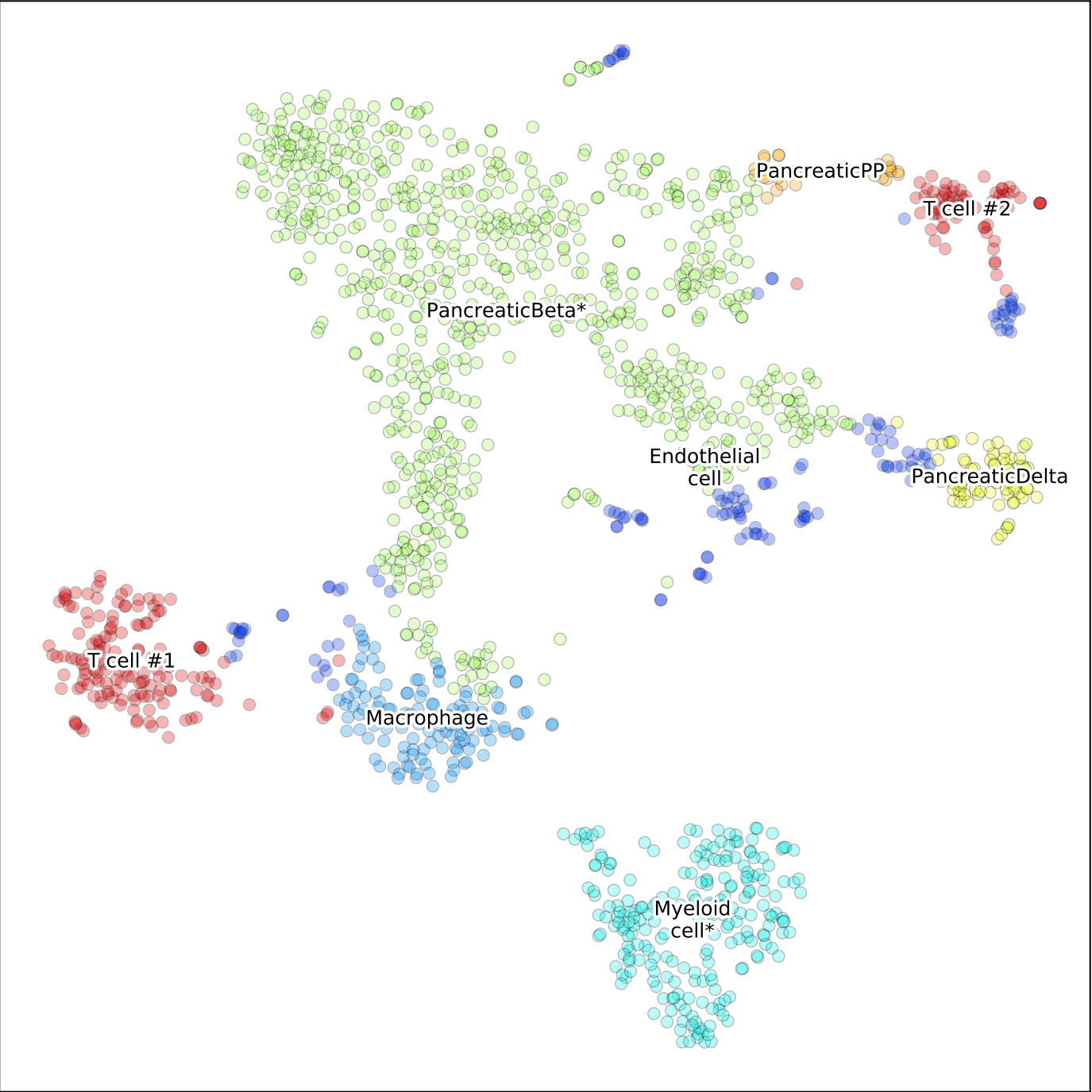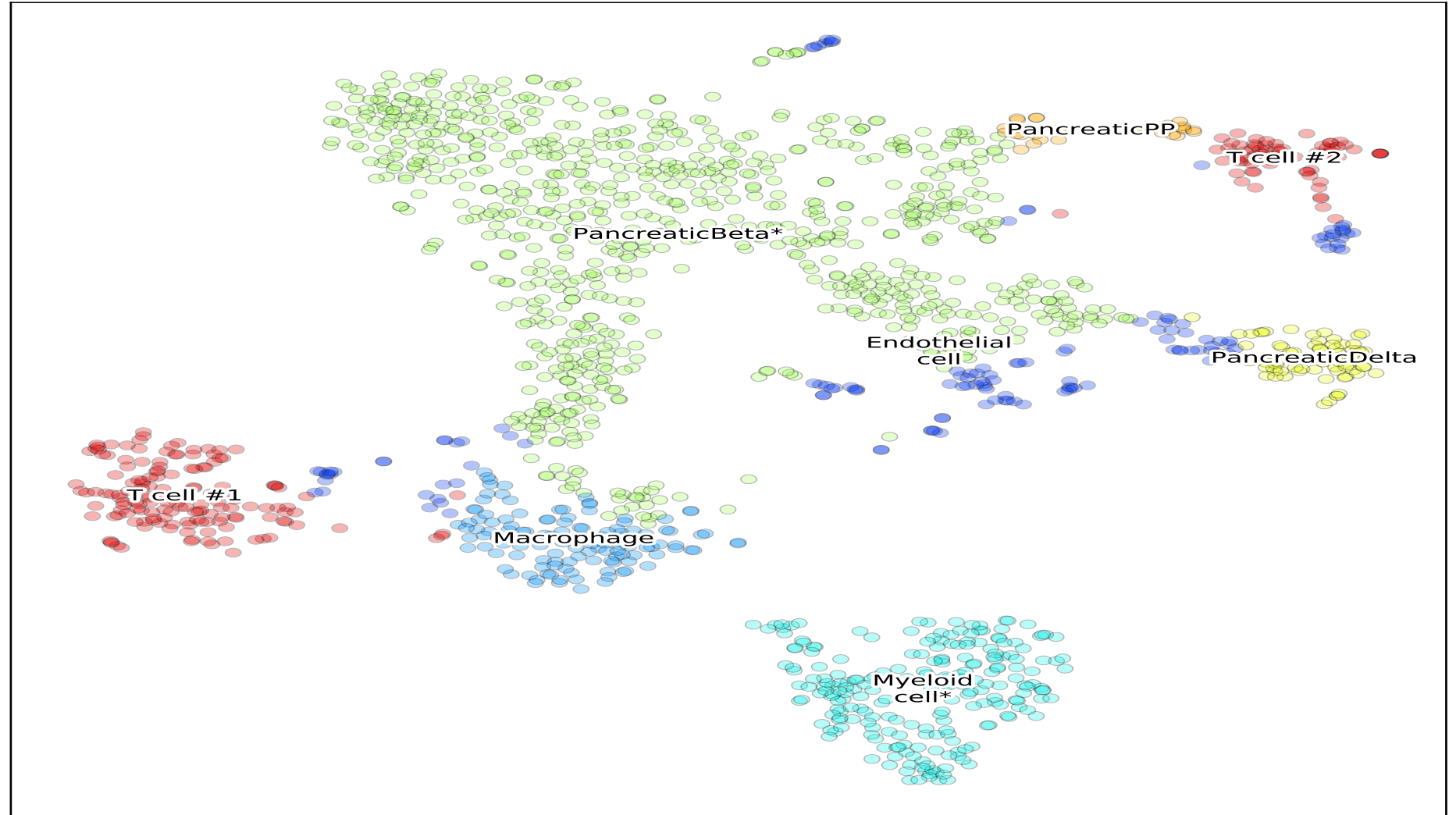| Labelled cell clusters in 2D/3D | Each maker expression in clusters |
| Voting results with cluster sizes | All marker/ce expression |

## Subcluster

To determine clusters' fine stru beginning for a desired cell type
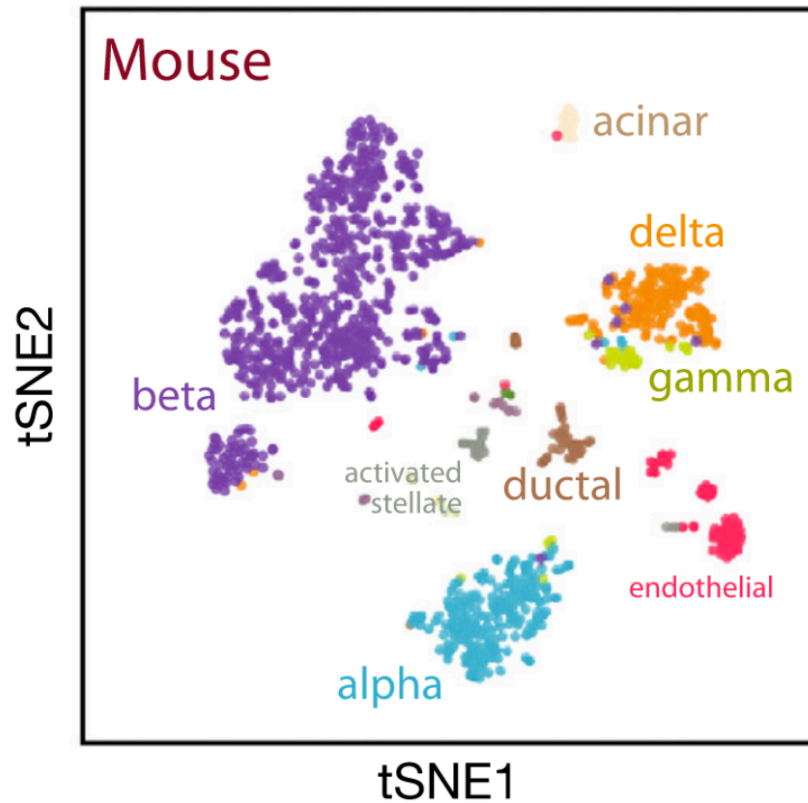
**DigitalCellSorter – Pancreas**

# scRNASeq – r package

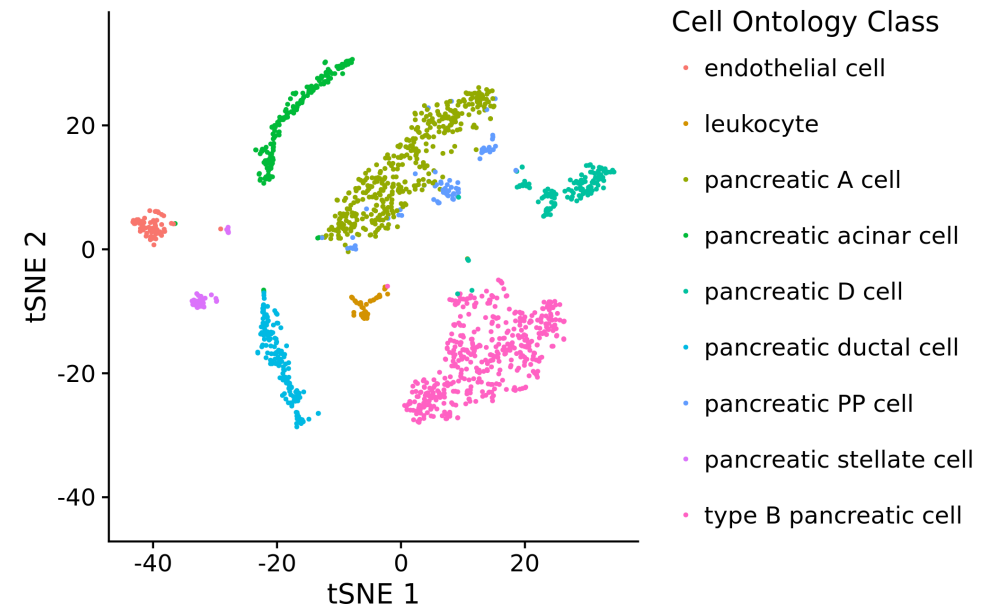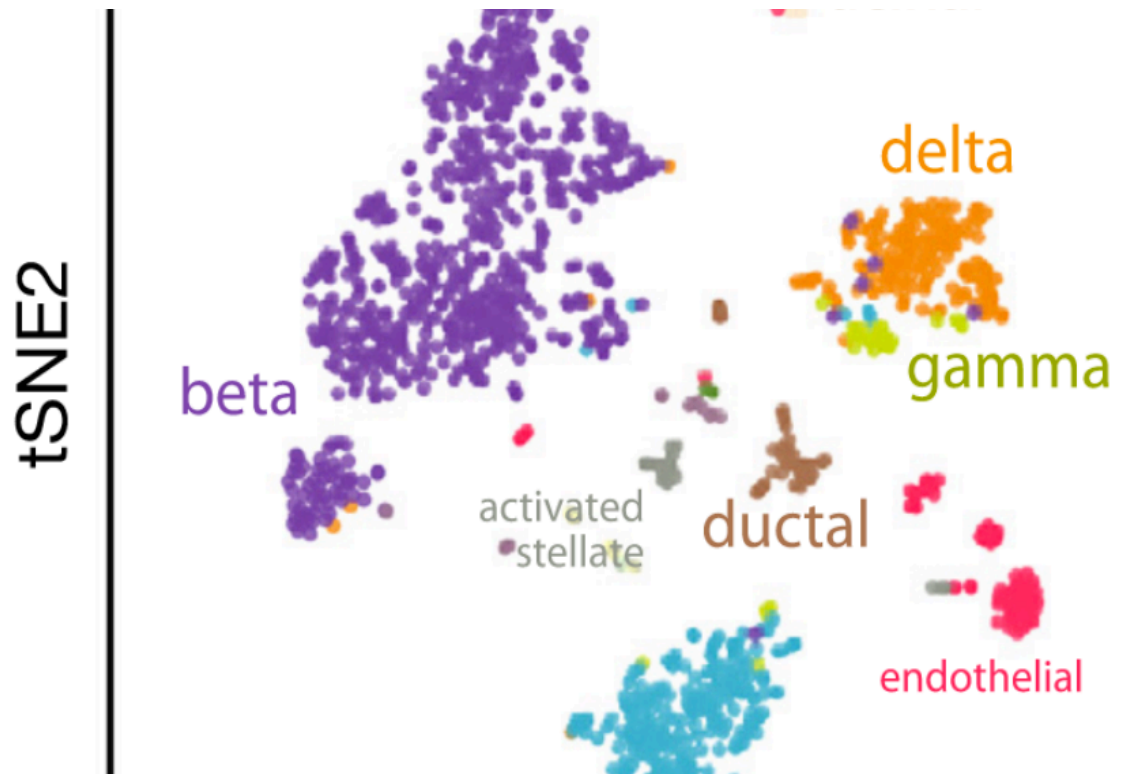| Study | Type | Cells |
|---|---|---|
| Lun et al. (2017) | 416B cells | 192 |
| La Manno et al. (2016) | Human embryonic midbrain | 1,977 |
| La Manno et al. (2016) | Human embryonic stem cells | 1,715 |
| Messmer et al. (2019) | Human embryonic stem cells | 1,344 |
| La Manno et al. (2016) | Human embyronic midbrain | 1,907 |
| La Manno et al. (2016) | Human induced pluripotent stem cells | 337 |
| Baron et al. (2016) | Human pancreas | 8,569 |
| Grun et al. (2016) | Human pancreas | 1,728 |
| Lawlor et al. (2017) | Human pancreas | 638 |
| Muraro et al. (2016) | Human pancreas | 3,072 |
| Segerstolpe et al. (2016) | Human pancreas | 3,514 |
| Xin et al. (2016) | Human pancreas | 1,600 |
| La Manno et al. (2016) | Mouse adult dopaminergic neurons | 243 |
| Campbell et al. (2017) | Mouse brain | 21,086 |
| Chen et al. (2017) | Mouse brain | 14,437 |
| Marques et al. (2016) | Mouse brain | 5,069 |
| Romanov et al. (2017) | Mouse brain | 2,881 |
| Usoskin et al. (2015) | Mouse brain | 864 |
| Tasic et al. (2016) | Mouse brain | 1,809 |
| Zeisel et al. (2015) | Mouse brain | 3,005 |
| Richard et al. (2018) | Mouse CD8+ T cells | 572 |
| Grun et al. (2016) | Mouse haematopoietic stem cells | 1,915 |
| Nestorowa et al. (2016) | Mouse haematopoietic stem cells | 1,920 |
| Bach et al. (2017) | Mouse mammary gland | 25,806 |
| Kolodziejczyk et al. (2015) | Mouse mebryonic stem cells | 704 |
| Baron et al. (2016) | Mouse pancreas | 1,886 |
| Macosko et al. (2015) | Mouse retina | 49,300 |
| Shekhar et al. (2016) | Mouse retina | 44,994 |
| Lun et al. (2017) | Mouse trophoblasts | 192 |
| Aztekin et al. (2019) | Xenopus tail | 13,199 |

# scRNASeq – r package



Baron et al. (2016)

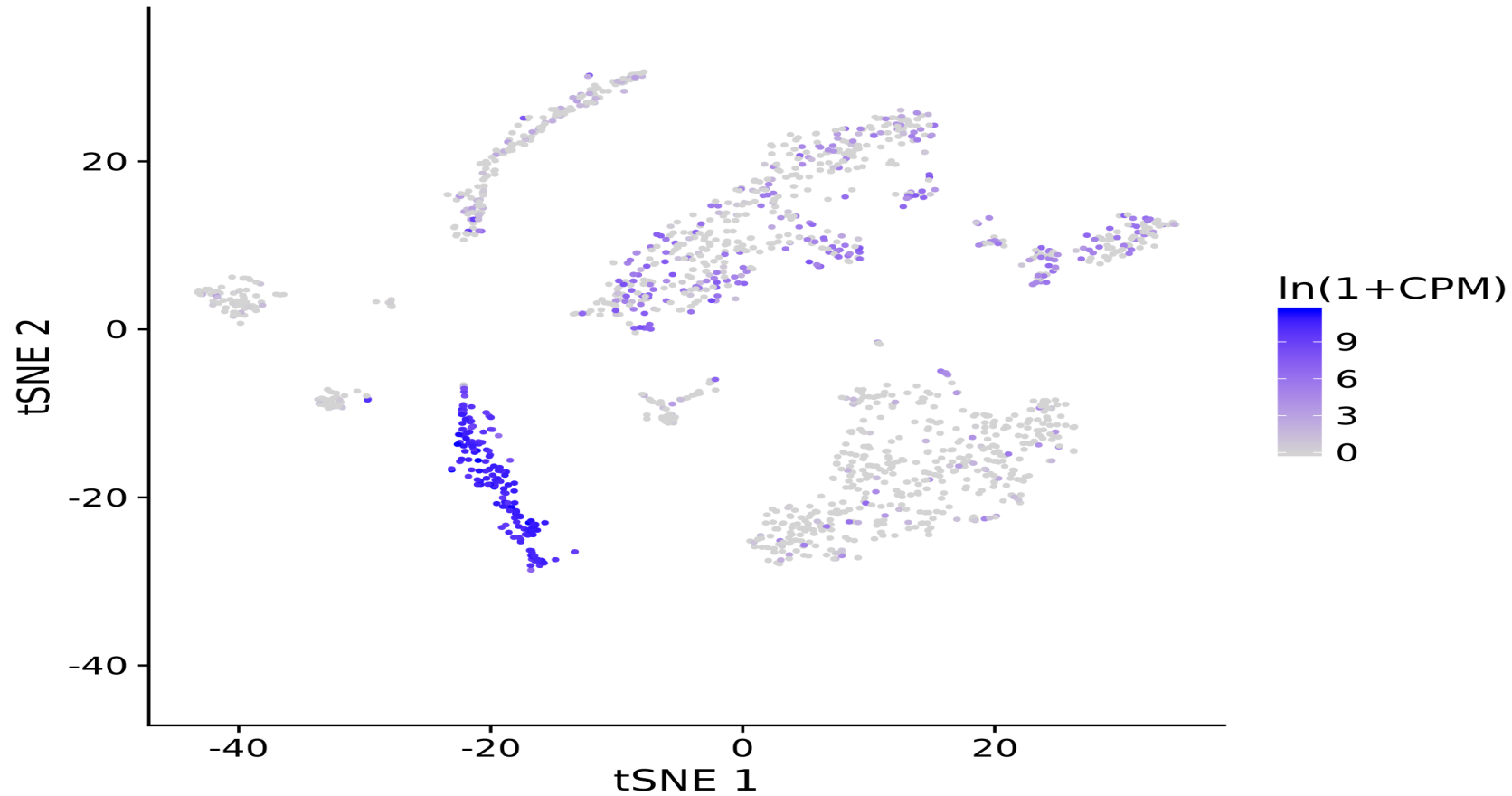| Study | Type | Cells |
|---|---|---|
| Lun et al. (2017) | 416B cells | 192 |
| La Manno et al. (2016) | Human embryonic midbrain | 1,977 |
| La Manno et al. (2016) | Human embryonic stem cells | 1,715 |
| Messmer et al. (2019) | Human embryonic stem cells | 1,344 |
| La Manno et al. (2016) | Human embyronic midbrain | 1,907 |
| La Manno et al. (2016) | Human induced pluripotent stem cells | 337 |
| Baron et al. (2016) - BM | Human pancreas | 8,569 |
| Grun et al. (2016) | Human pancreas | 1,728 |
| Lawlor et al. (2017) | Human pancreas | 638 |
| Muraro et al. (2016) | Human pancreas | 3,072 |
| Segerstolpe et al. (2016) | Human pancreas | 3,514 |
| Xin et al. (2016) | Human pancreas | 1,600 |
| La Manno et al. (2016) | Mouse adult dopaminergic neurons | 243 |
| Campbell et al. (2017) | Mouse brain | 21,086 |
| Chen et al. (2017) | Mouse brain | 14,437 |
| Marques et al. (2016) | Mouse brain | 5,069 |
| Romanov et al. (2017) | Mouse brain | 2,881 |
| Usoskin et al. (2015) | Mouse brain | 864 |
| Tasic et al. (2016) | Mouse brain | 1,809 |
| Zeisel et al. (2015) | Mouse brain | 3,005 |
| Richard et al. (2018) | Mouse CD8+ T cells | 572 |
| Grun et al. (2016) | Mouse haematopoietic stem cells | 1,915 |
| Nestorowa et al. (2016) | Mouse haematopoietic stem cells | 1,920 |
| Bach et al. (2017) | Mouse mammary gland | 25,806 |
| Kolodziejczyk et al. (2015) | Mouse mebryonic stem cells | 704 |
| Baron et al. (2016) | Mouse pancreas | 1,886 |
| Macosko et al. (2015) | Mouse retina | 49,300 |
| Shekhar et al. (2016) | Mouse retina | 44,994 |
| Lun et al. (2017) | Mouse trophoblasts | 192 |
| Aztekin et al. (2019) | Xenopus tail | 13,199 |

# Tabula Muris Consortium (TM)

The **Chan Zuckerberg Biohub** recently released **Tabula Muris**, a compendium of single cell transcriptome data from the mouse containing nearly **100,000 cells from 20 organs and tissues**. The data allow for direct and controlled comparison of gene expression in cell types shared between tissues, such as immune cells from distinct anatomical locations. They also allow for a comparison of two distinct technical approaches:
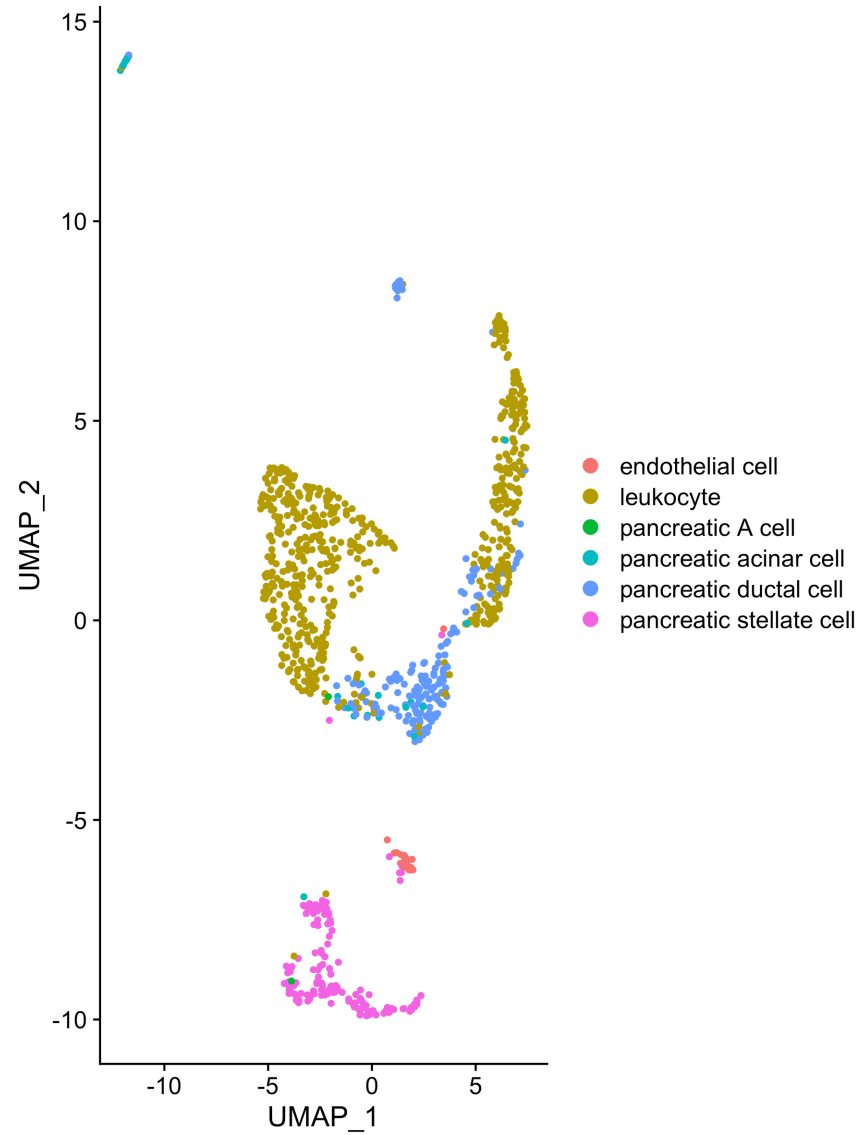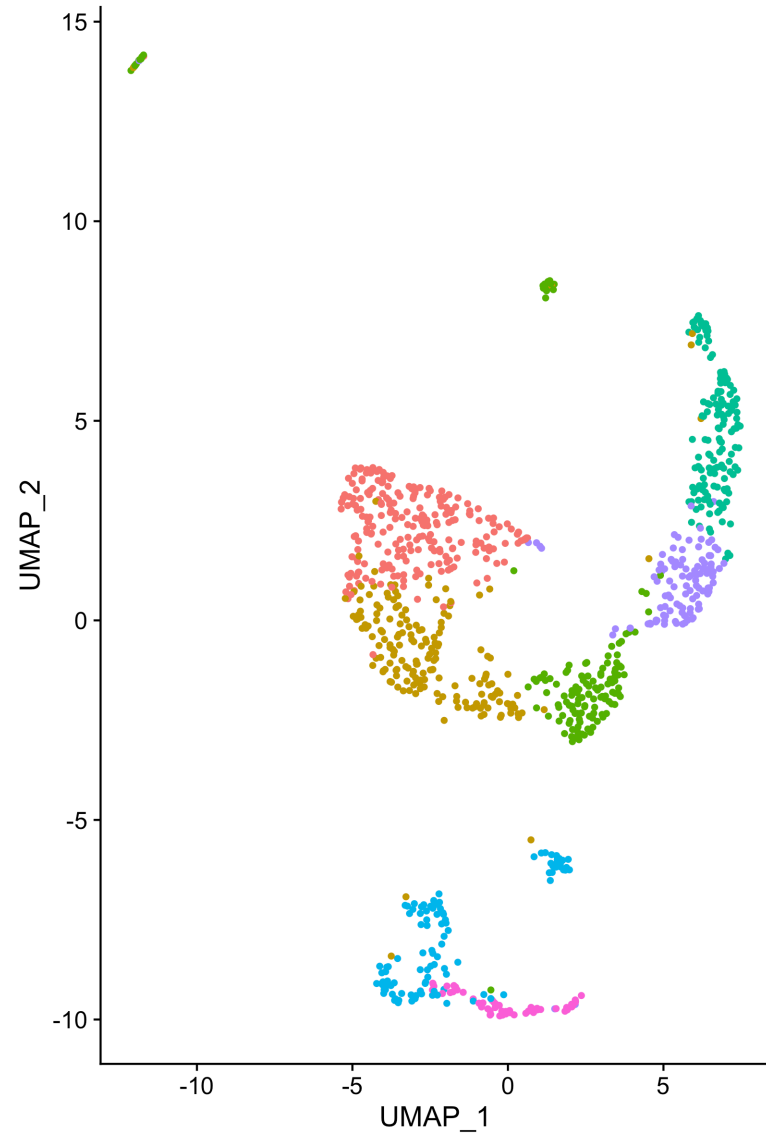
- microfluidic droplet-based 3'-end counting, which provides a survey of thousands of cells per organ at relatively low coverage.
- FACS-based full length transcript analysis, which provides higher sensitivity and coverage.

# SingleR – Label Transfer - Tabula Muris (TM)

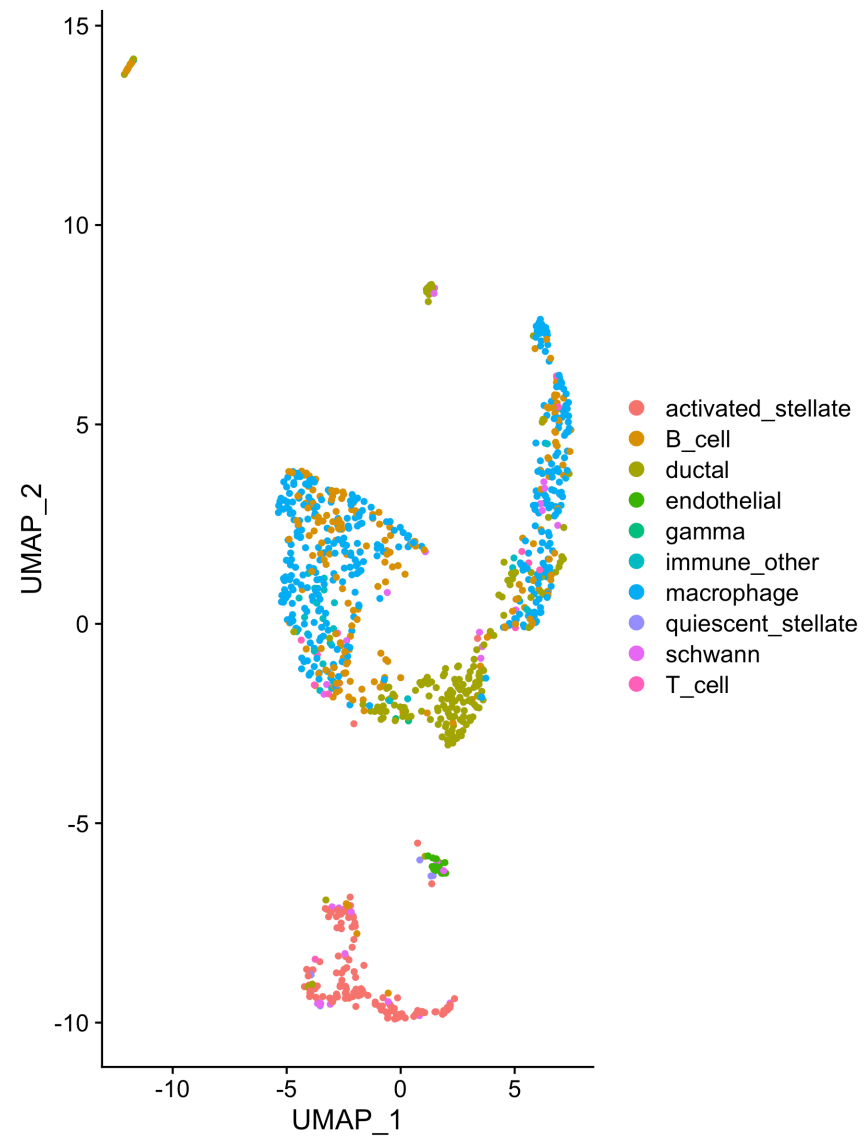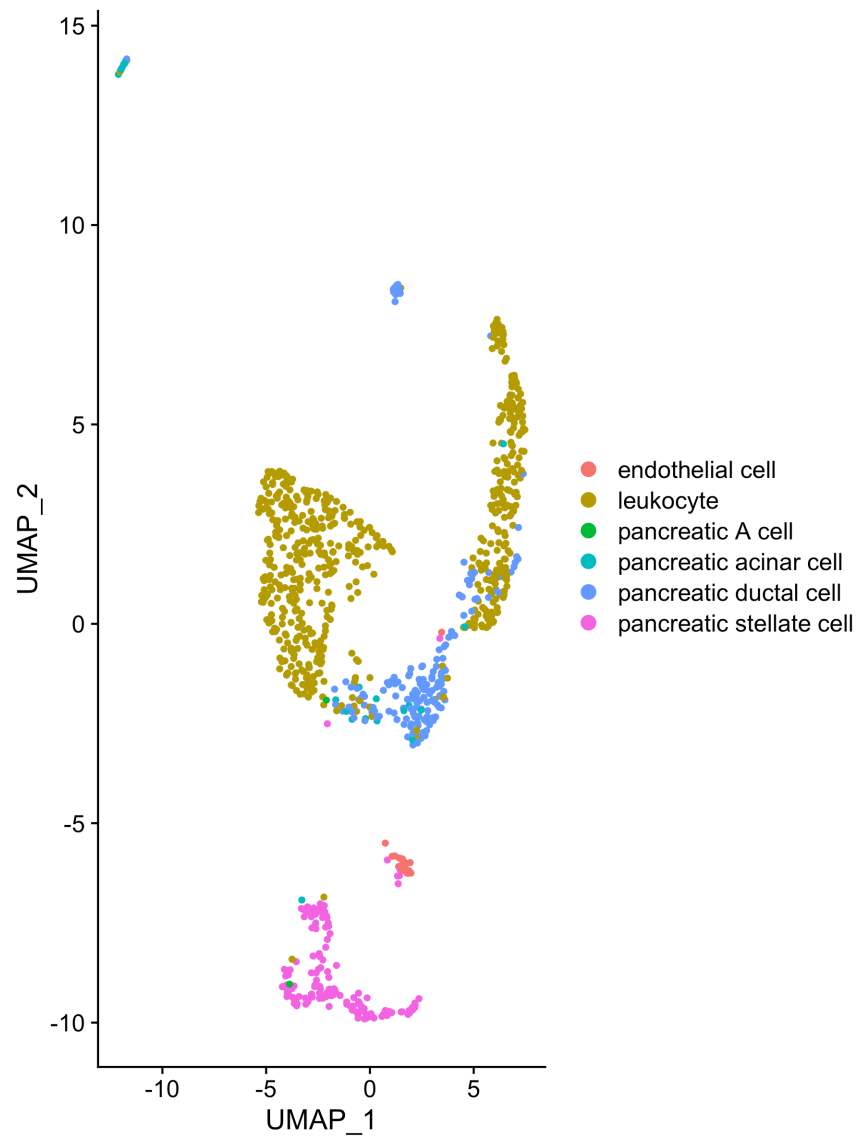# SingleR – Label Transfer – Barron et al (BM)

# SingleR – Label Transfer – BM vs TM

# SingleR – Label Transfer – BM vs TM
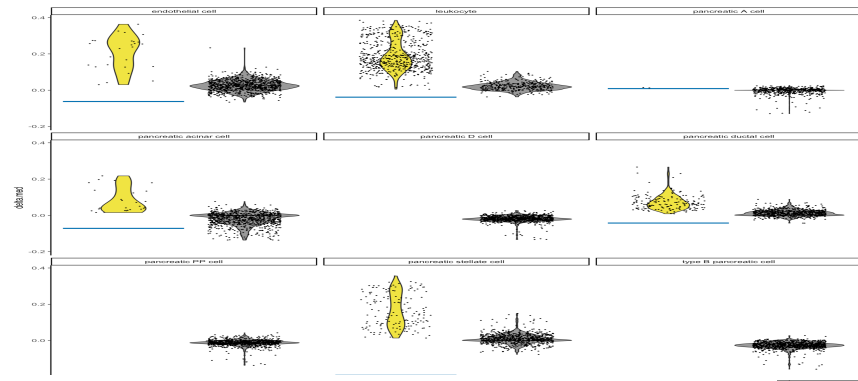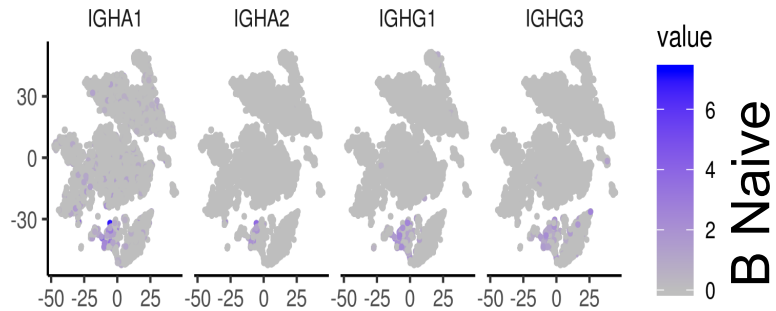
# SF scRNA Pipeline Automated Panels

- B Pan
- B Memory
- B Naive
- Macrophage
- NK
- T Central Memory
- T Cytotoxic
- T Effector Memory
- NKT
- T Pan
- T Regulatory
- T Tfh
- T Th1

- T Th1
- T Th17
- T Th2
- T gd
- Monocyte Ly-6C hi
- Monocyte Ly-6C lo
- Cell Cycle M
- Cell Cycle S
- Pancreas
- Pancreatic Alpha
- Pancreatic Beta
- Pancreatic Delta
- Pancreatic PP
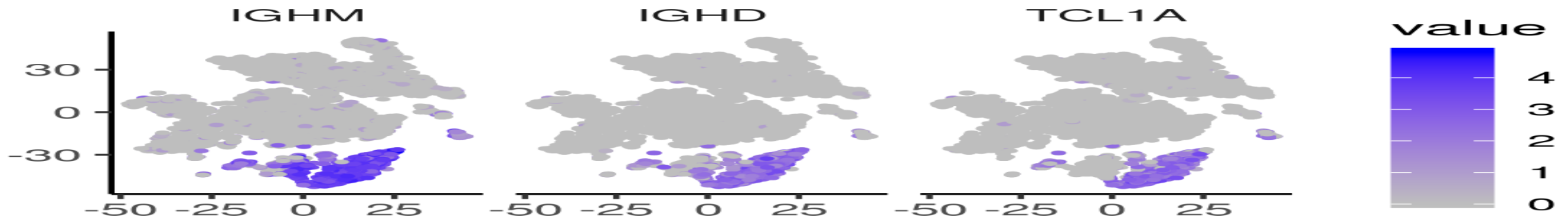- Pancreatic Duct

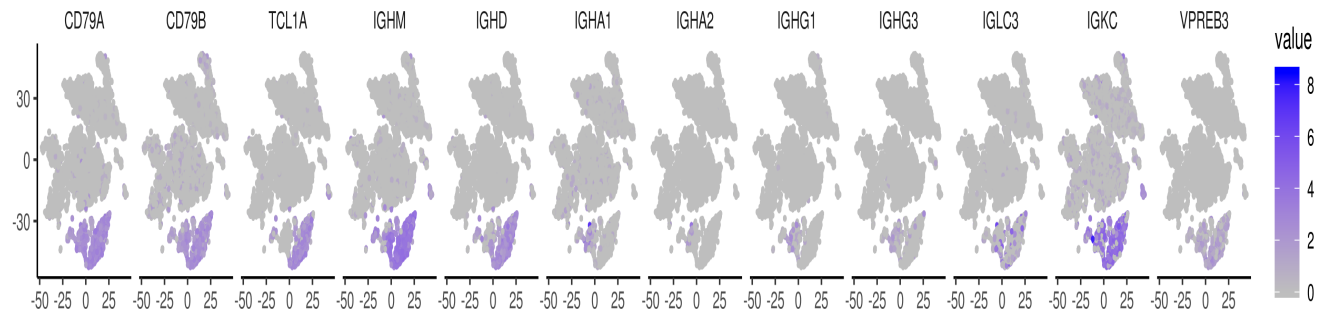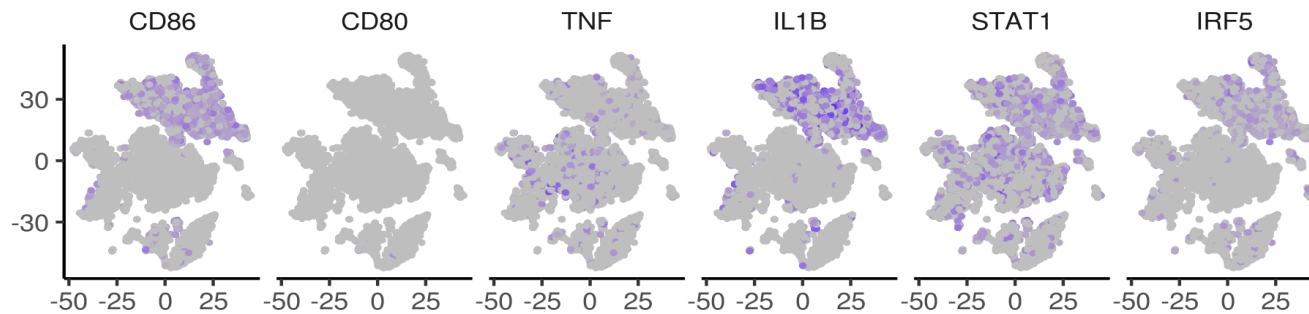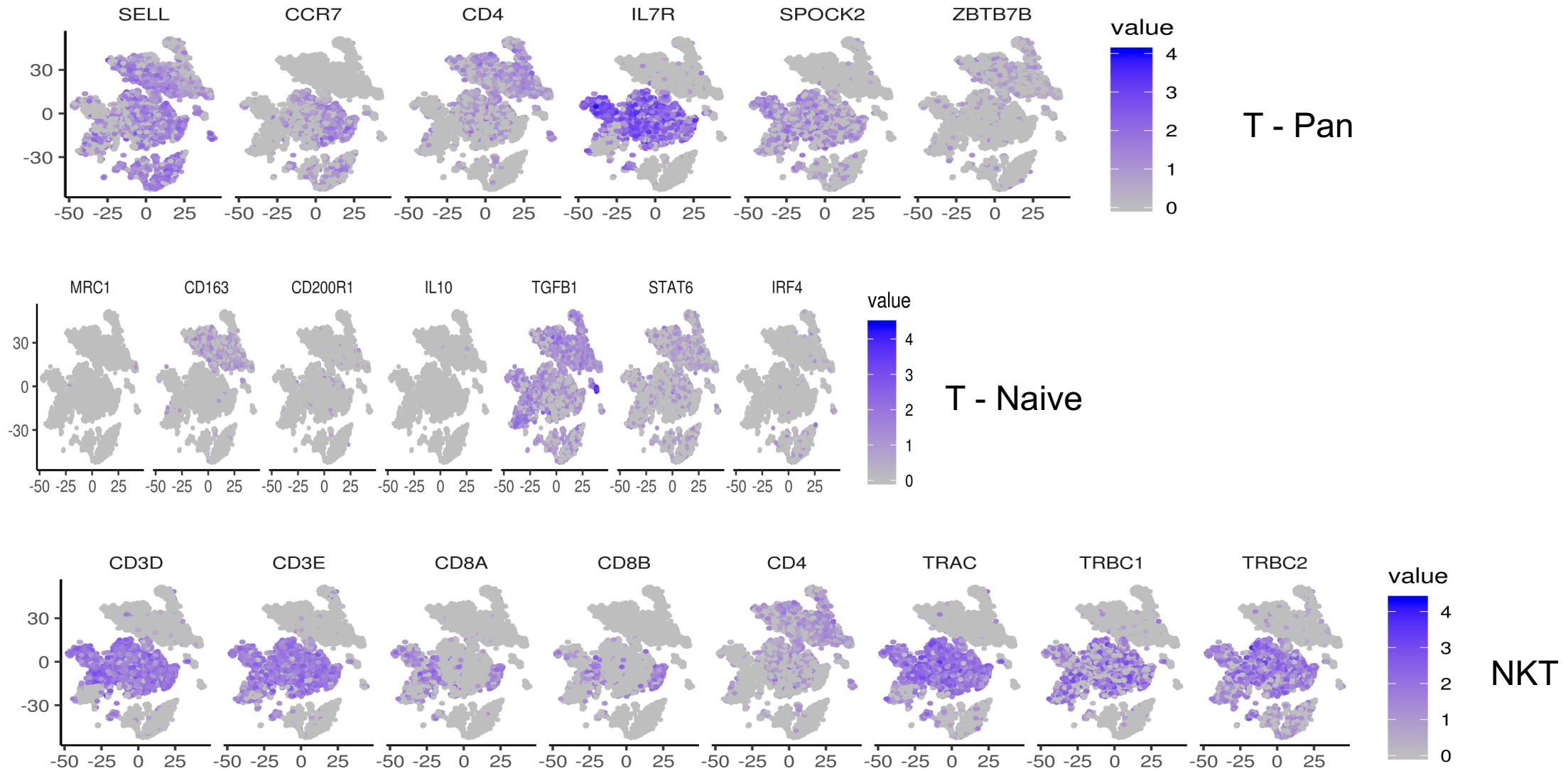# SF scRNA Pipeline Automated Panels



Macrophage – M1

Macrophage – M2

# SF scRNA Pipeline Automated Panels

# Question?