# Single-cell analysis workshop

## Sydney Precision Bioinformatics Group

Workshop presenters:
Hani Kim
Yingxin Lin
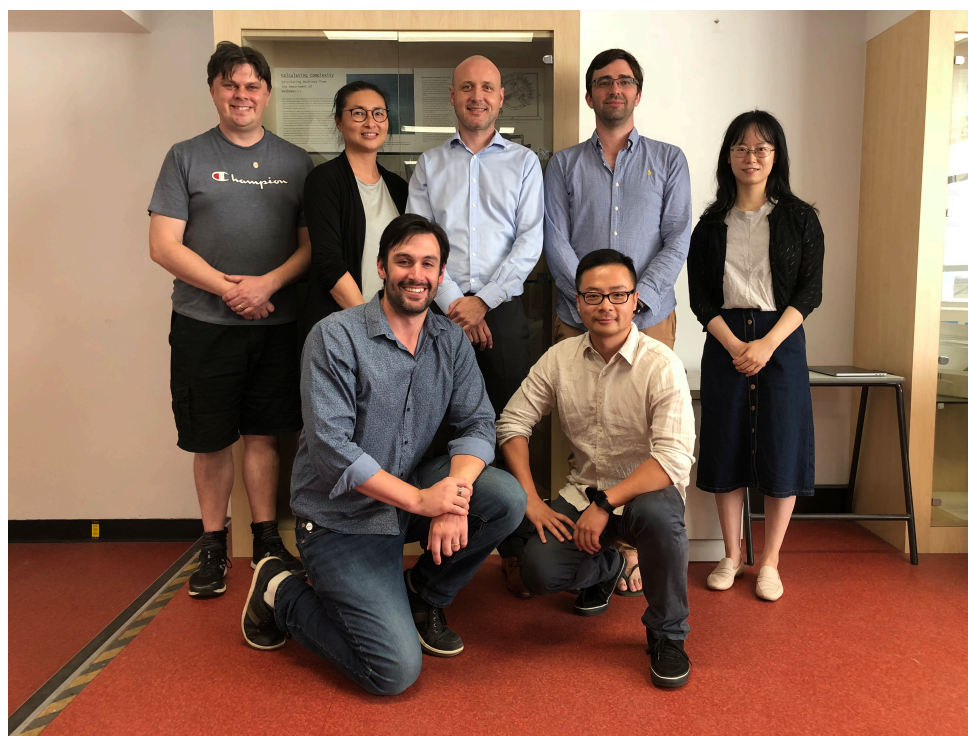Shila Ghazanfar

THE UNIVERSITY OF
SYDNEY

# Sydney Precision Bioinformatics Research Group

We share an interest in developing statistical and computational methodologies to tackle the foremost significant challenges posed by modern biology and medicine.

Meet our senior and junior research leaders:

A/Prof. John Ormerod; Prof. Jean Yang; Prof. Samuel Mueller; Dr. Garth Tarr; Dr. Rachel Wang



Dr. Ellis Patrick; Dr. Pengyi Yang

and senior research associates, PhD candidates, Honours and TSP students.

Find out more:

http://www.maths.usyd.edu.au/bioinformatics/

Shiny apps:     http://shiny.maths.usyd.edu.au/

Github:     https://github.com/SydneyBioX

# Roadmap for the workshop

Setting up: 13:30 – 13:45 Google cloud set up

Session 1: 13:45 – 14:15 Single cell analysis overview (scdney)

Session 2: 14:15 – 15:00 Quality control and data integration

AFTERNOON TEA: 1500-1530

Session 3: 15:30 – 16:00 Overview of single-cell downstream analysis

Session 4: 16:00 – 16:45 Downstream analysis: cell type identification, identify marker genes & cell type composition

Extension: cell type identification via supervised classification and single cell trajectory analysis

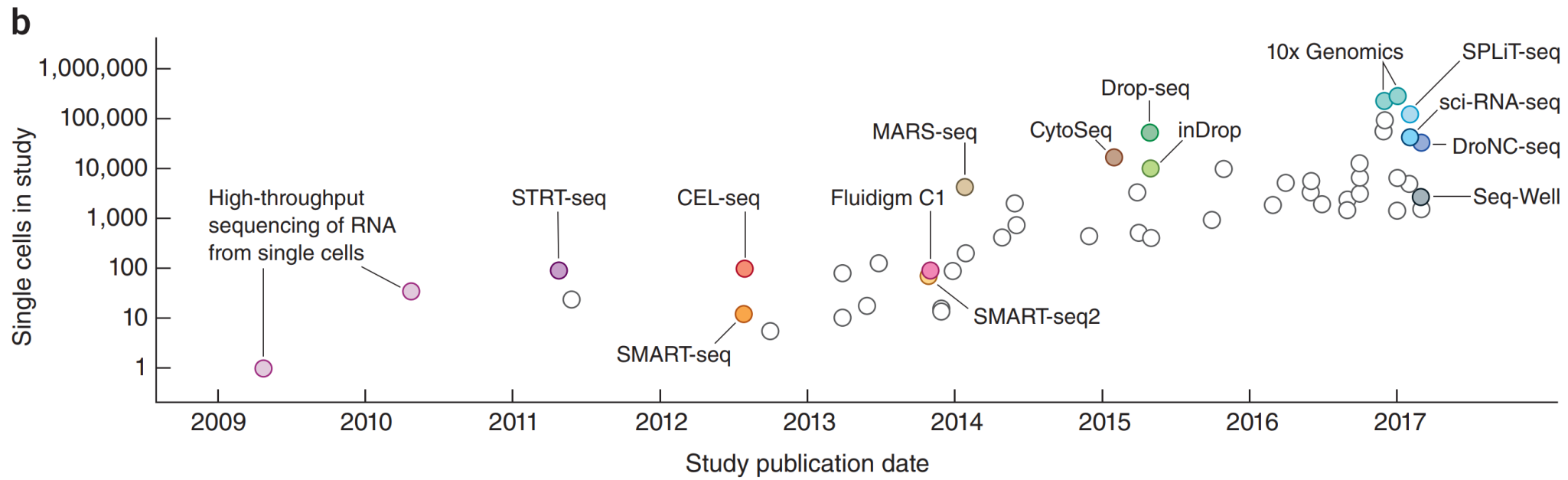**Configuring Google Cloud and workshop materials**
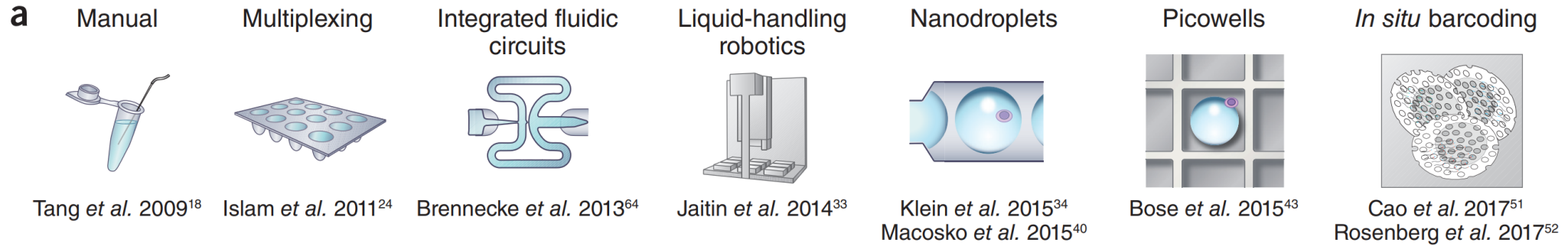
– Workshop materials:
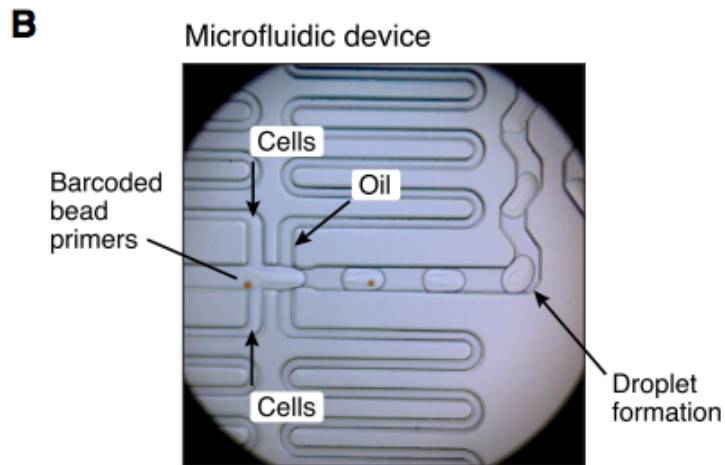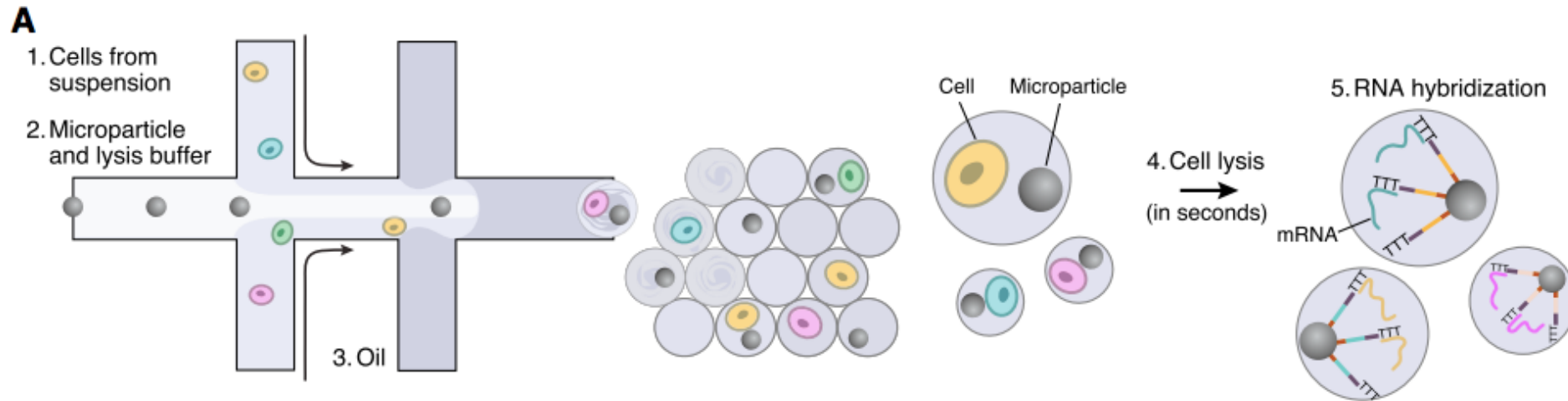https://sydneybiox.github.io/BIS2019_SC/index.html

–Machine 1: 34.68.240.36
–Machine 2: 34.94.37.174
source("/home/user_setup.R")

# Exponential growth in single cell RNA seq technologies



**a**

| Manual | Multiplexing | Integrated fluidic circuits | Liquid-handling robotics | Nanodroplets | Picowells | *In situ* barcoding |
|---|---|---|---|---|---|---|
| Tang *et al.* 2009[18] | Islam *et al.* 2011[24] | Brennecke *et al.* 2013[64] | Jaitin *et al.* 2014[33] | Klein *et al.* 2015[34] Macosko *et al.* 2015[40] | Bose *et al.* 2015[43] | Cao *et al.* 2017[51] Rosenberg *et al.* 2017[52] |

Svensson et al. *Nature Protocols* (2018)

# Droplet based technologies are now dominating



Macosko et al. (2015), *Cell*

10X Genomics is a commercial provider of droplet based scRNAseq platform

# scRNAseq experiments approaching 1 million cells



Saunders et al., (2018) Cell

**690,000 individual cells** from 9 regions
of adult mouse brain

# Number of scRNAseq tools also increasing rapidly
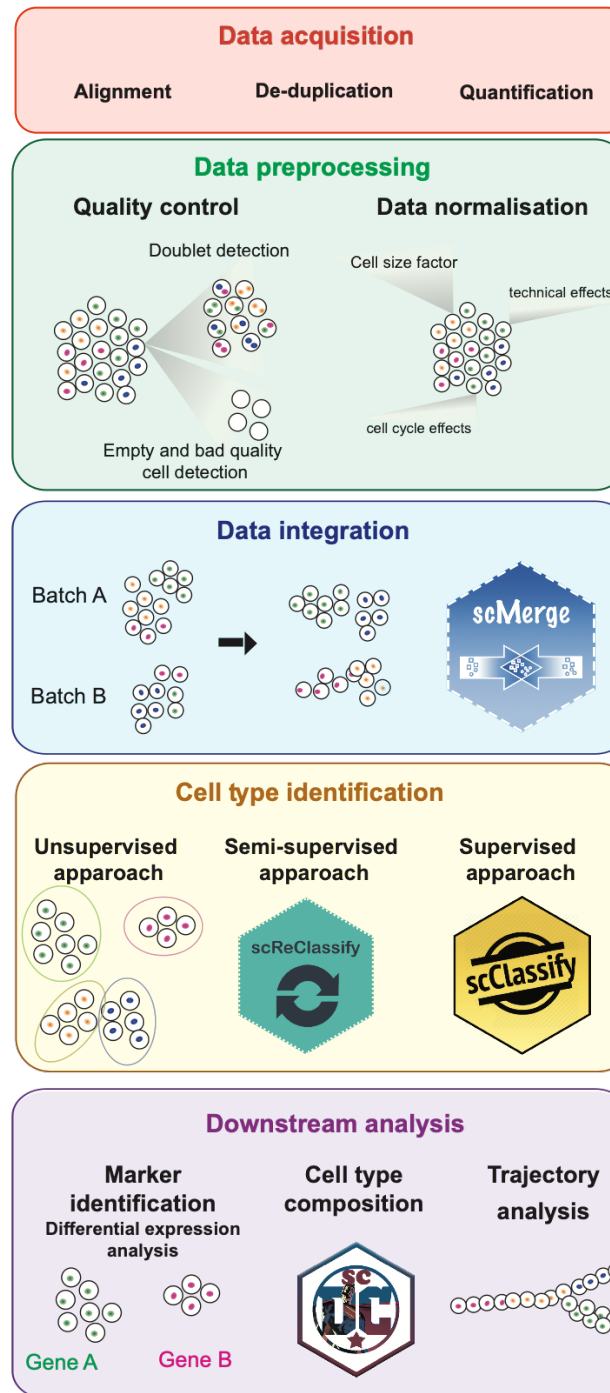


Downloaded from www.scrna-tools.org

# Single-cell RNA-seq analysis

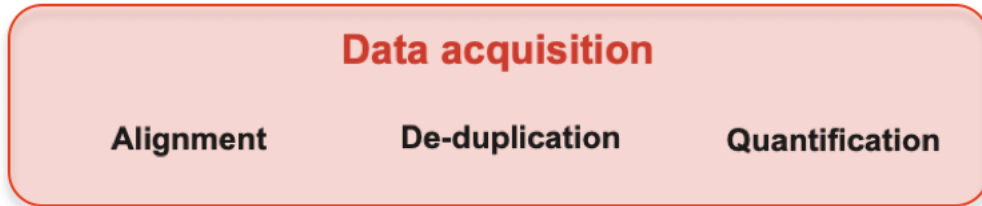# Components of a typical scRNA-seq analysis process

# Component 1: Data acquisition



Data acquisition
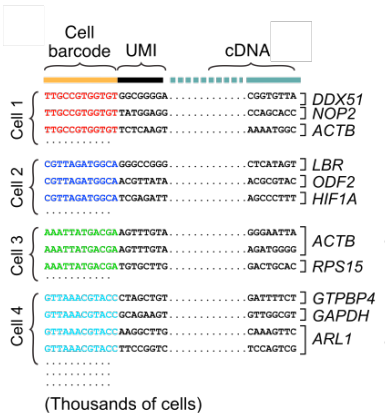
Alignment    De-duplication    Quantification

## Input
- BCL or fastq file from the sequencer

## Output
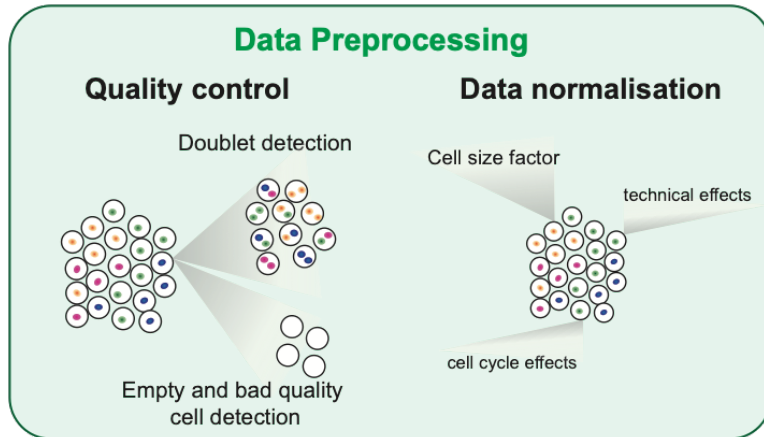- Gene/cell counts matrix



(Thousands of cells)

## Software
- CellRanger for 10X Genomics data
- Macosko's custom scripts for DropSeq data
- STAR for alignment plus custom scripts (or there is STAR-solo)

## Considerations
- Single or mix of species? Does it include ERCC spike-ins? May need to build a custom reference
- Barcode and/or UMI sequencing errors – CellRanger takes care of this automatically
- Align to exon or exon and intron?

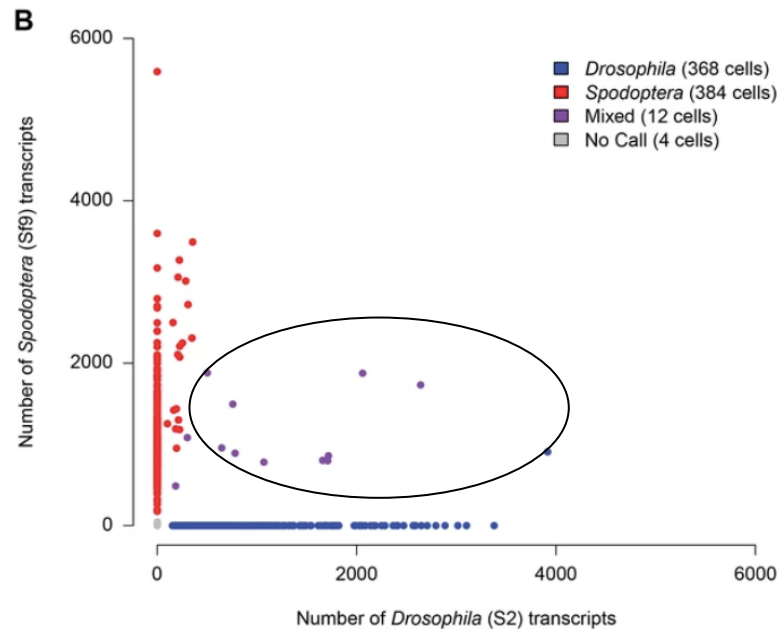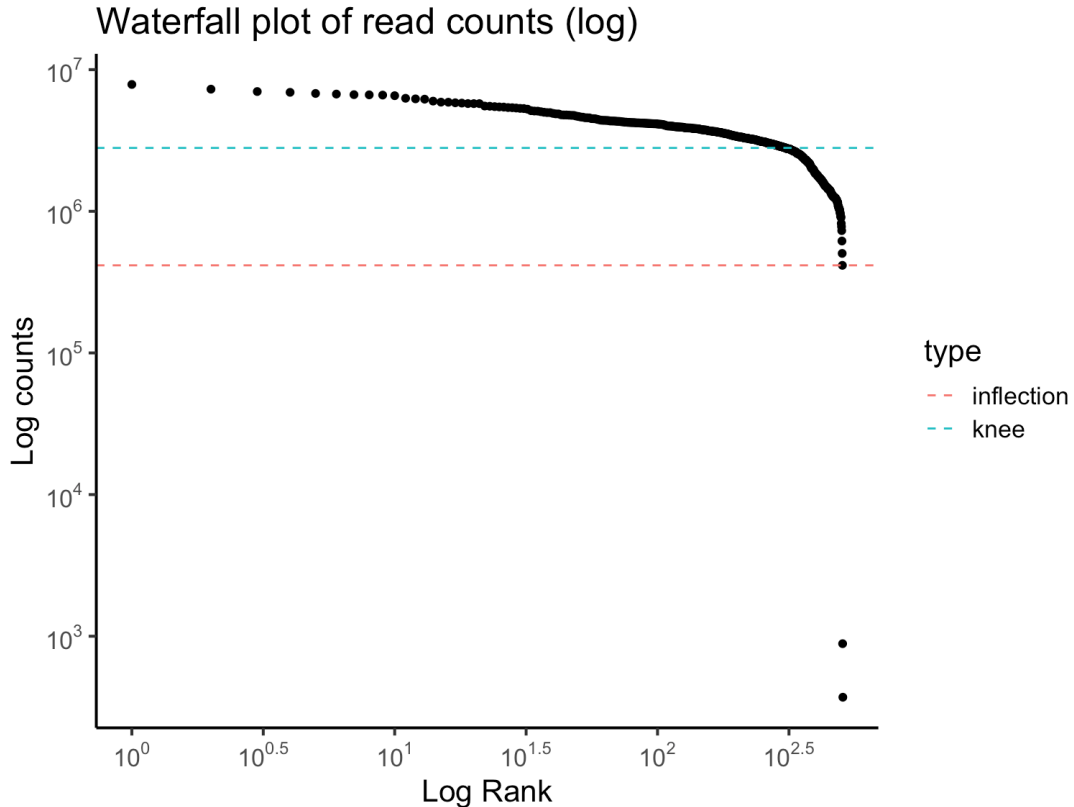# Component 2: Data preprocessing – Quality control



## Software
- Seurat (all-purpose single cell R package)
- Scater
- DropletUtils (R package with a number of handy utility functions)
- Your own custom scripts

## Considerations
- Filter out droplets with doublets – may be difficult to find. Can estimate expected rate by doing species mixture experiment

Croset (2018), eLife

# Component 2: Data preprocessing – Quality control
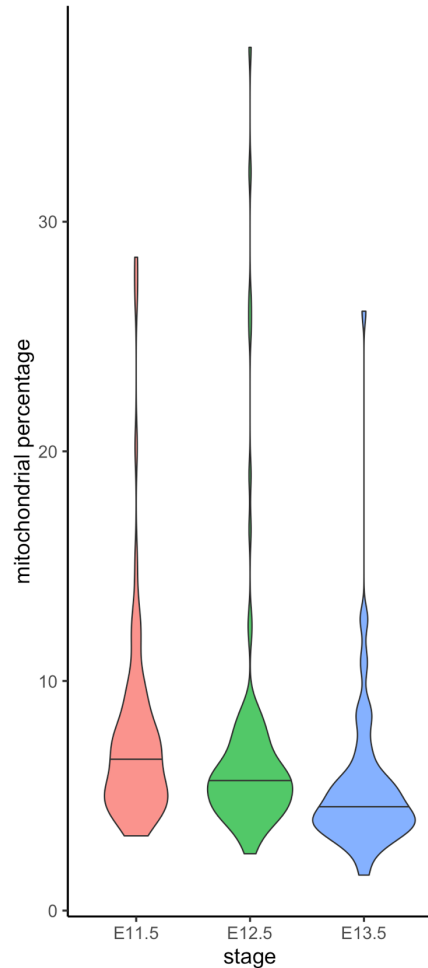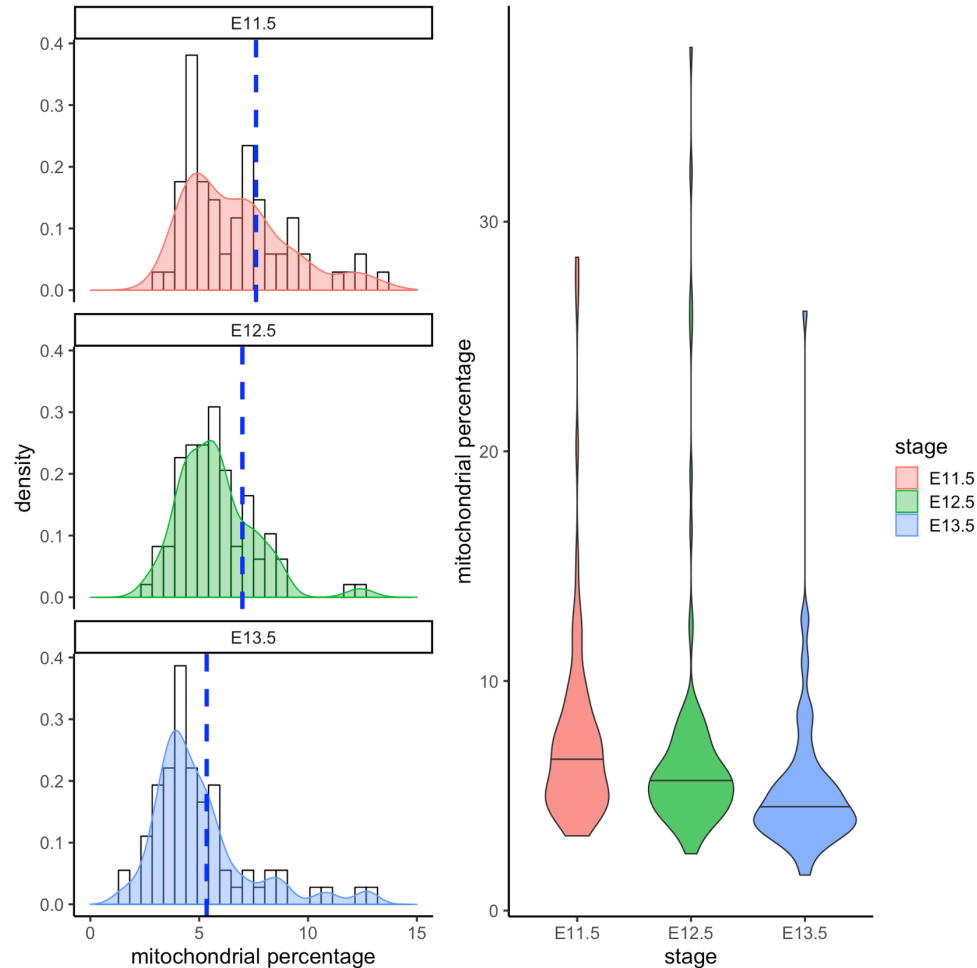
## Waterfall plot of read counts (log)



## Software
- Seurat (all-purpose single cell R package)
- Scater
- DropletUtils (R package with a number of handy utility functions)
- Your own custom scripts

## Considerations
- Filter out droplets with doublets – may be difficult to find. Can estimate expected rate by doing species mixture experiment
- Filter out droplets with no cells

# Component 2: Data preprocessing – Quality control



## Software

- Seurat (all-purpose single cell R package)
- Scater
- DropletUtils (R package with a number of handy utility functions)
- Your own custom scripts

## Considerations

- Filter out droplets with doublets – may be difficult to find. Can estimate expected rate by doing species mixture experiment
- Filter out droplets with no cells
- Filter out droplets with damaged cells – look for high mitochondrial gene content or high spike-in

# Component 2: Data normalisation

## Software
- scran for non-full-length datasets (Lun et al. Genome Biology 2016)
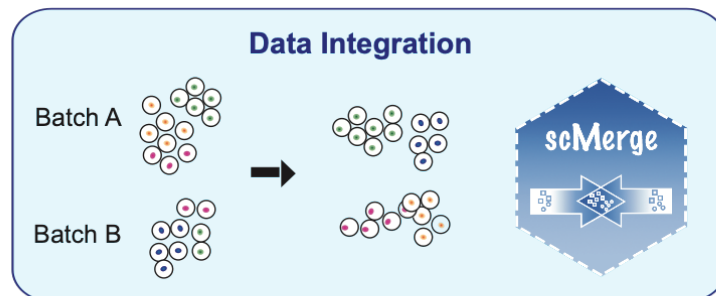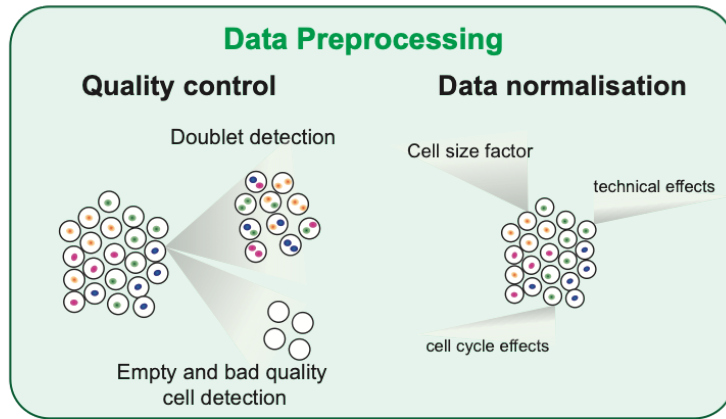- bulk methods for full-length datasets (TPM normalisation)

## Normalisation aims to address
- Removing sampling effects
- Scaling count data to obtain correct relative gene expression abundances between cells

## After normalisation, data matrices are typically log(x+1)-transformed
- Distances represent log-fold changes
- log transformation mitigates (but does not remove) the mean–variance relationship in single-cell data
- reduces the skewness of the data

# Component 3: Data integration



**Data Preprocessing**

**Quality control**

Doublet detection

Empty and bad quality cell detection

**Data normalisation**

Cell size factor

technical effects

cell cycle effects



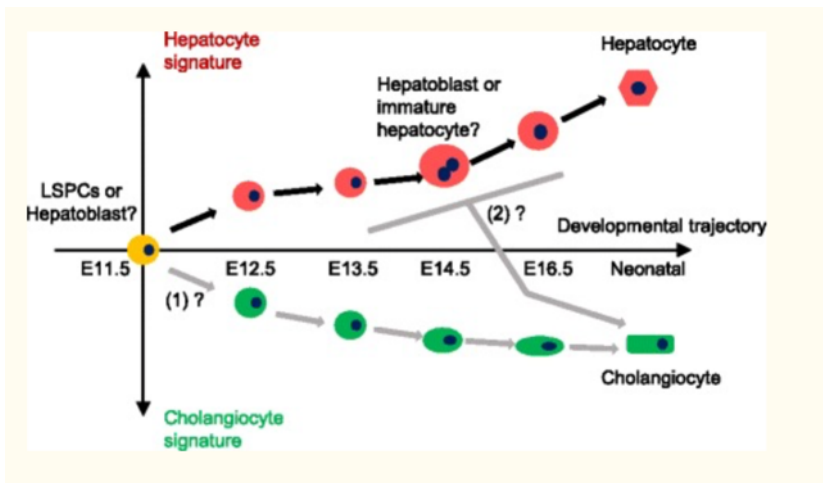**Data Integration**

Batch A

Batch B

scMerge

**Software**

- Seurat (all-purpose single cell R package) for very basic normalization
- Batch effect correction
  - mnnCorrect
  - Harmony
  - Liger
  - **scMerge**

# Liver fetal development time course datasets

# tSNE of liver fetal development time course datasets



**Highlighted by cell types**

**Highlighted by batches**

Challenge:
Strong "batch effect"

Cell types: cholangiocyte, Epithelial Cell, hepatoblast/hepatocyte, Mesenchymal Cell, Endothelial Cell, Hematopoietic, Immune cell, Stellate Cell

Batch: GSE87038, GSE87795, GSE90047, GSE96981

# Breaking observed data into components

For *n* cells with data collected for *m* genes

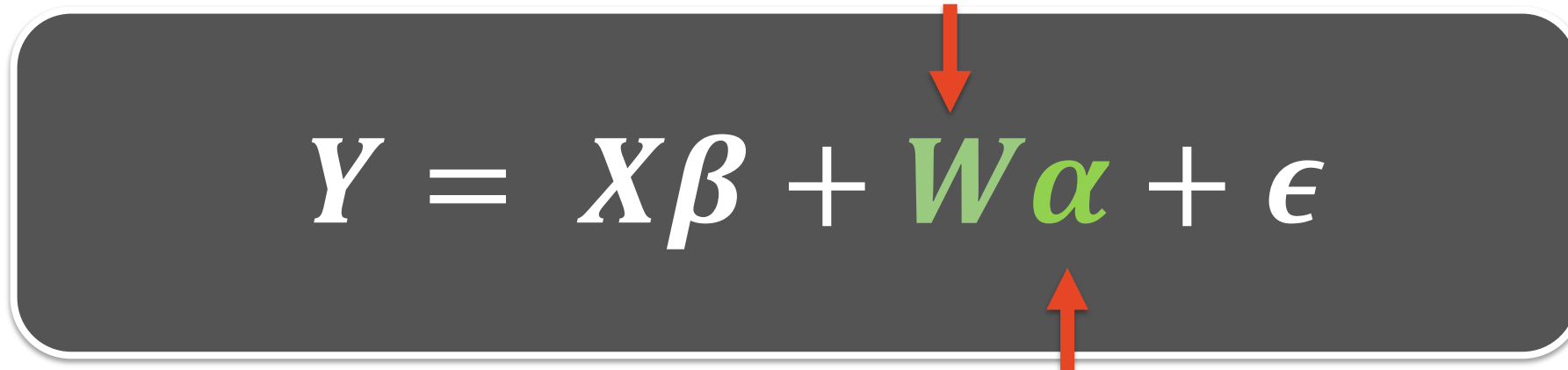$$Y = X\beta + W\alpha + \epsilon$$

The data we observe

Biologically relevant variation
cell types
p *wanted* variables

Unwanted variation batch and technical effects
k *unwanted* variables

Random noise

# scMerge algorithm

Estimated by **stably expressed genes** by factor analysis
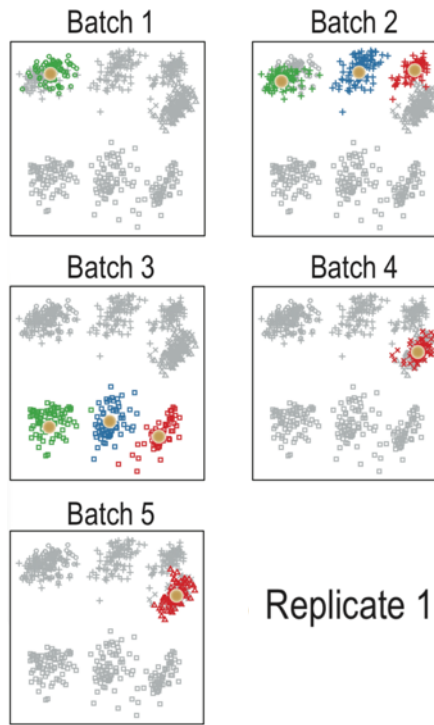
$$Y = X\beta + W\alpha + \epsilon$$

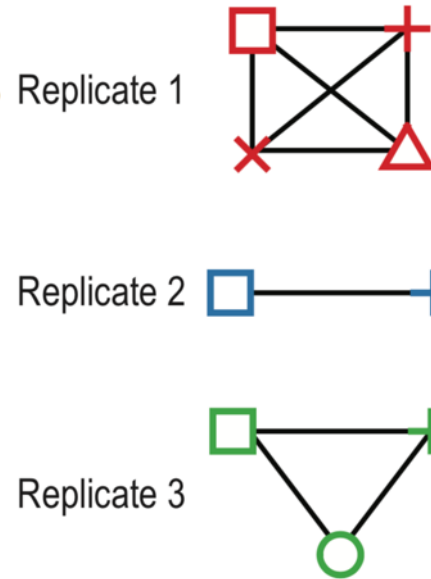Estimated with **replicates** by factor analysis

RUVIII algorithm Molania et al. (2019), Nuclei Acids Res

# scMerge algorithm



**Clustering** for each batch
(k-means by default)
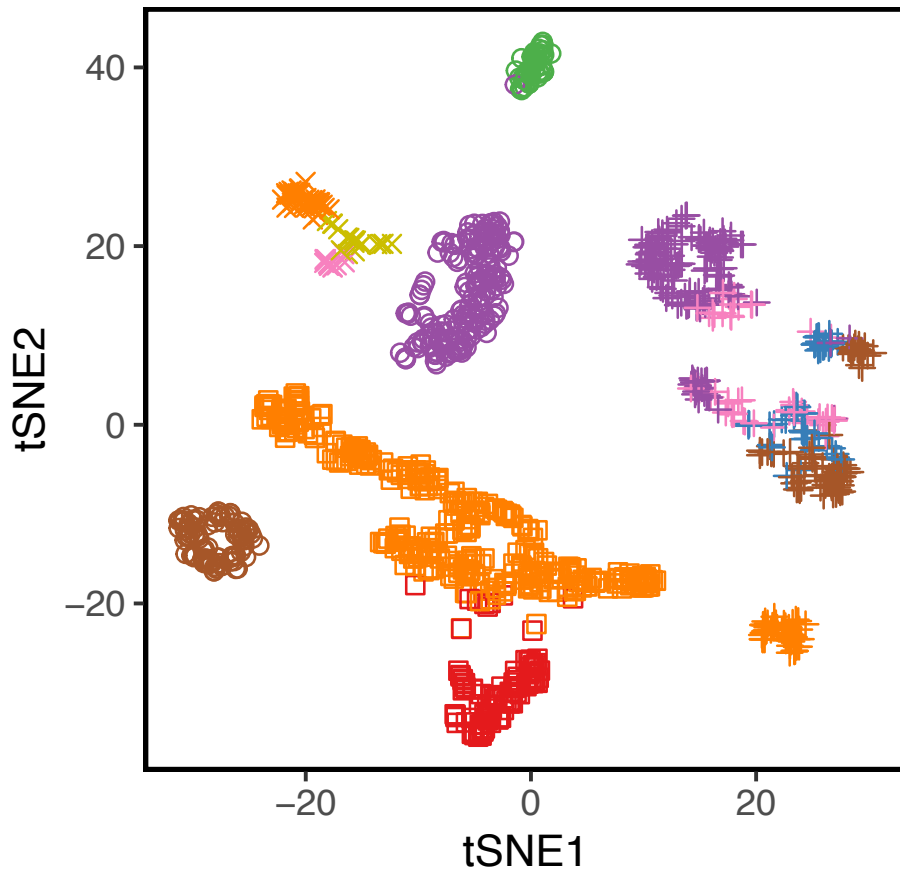
Find **Mutual Nearest Clusters** as pseudo-replicates
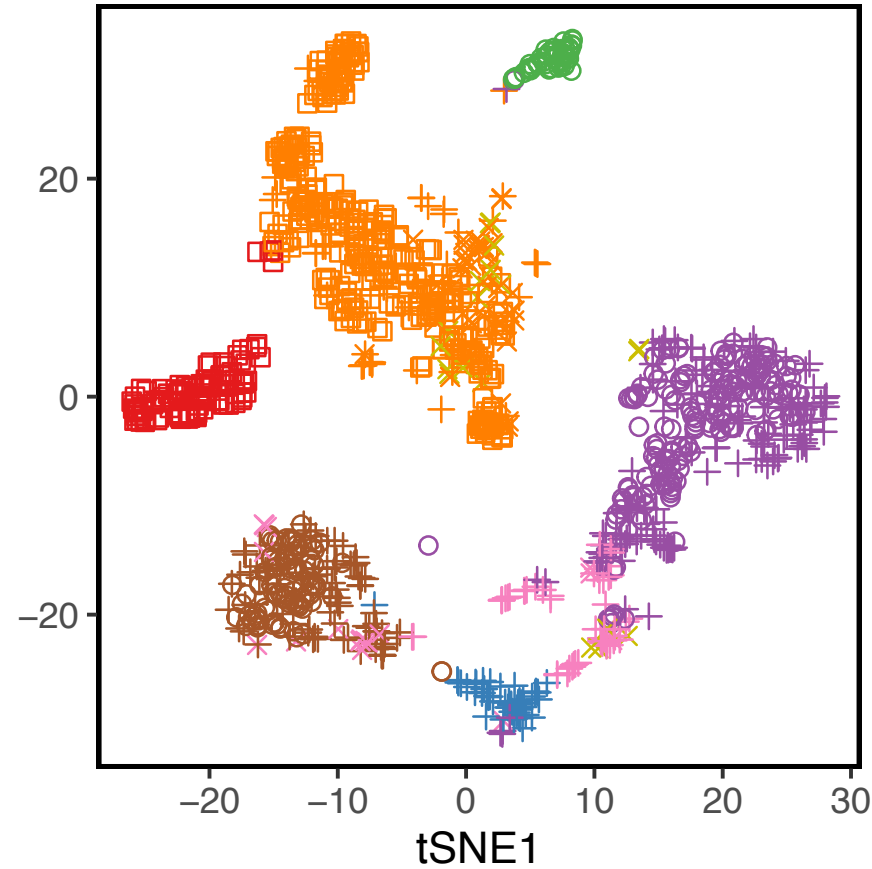
**Frame as pseudo-replicate information**

Pseudo-replicates

# Coming back to our motivational data –
# Liver fetal development time course datasets



Before scMerge / After scMerge

cell_types
- cholangiocyte
- Endothelial Cell
- Epithelial Cell
- Hematopoietic
- hepatoblast/hepatocyte
- Immune cell
- Mesenchymal Cell
- Stellate Cell

batch
- GSE87038
- GSE87795
- GSE90047
- GSE96981

# More information

**PNAS:**
https://doi.org/10.1073/pnas.1820006116

**scMerge R package and website:**
https://sydneybiox.github.io/scMerge/

# SingleCellExperiment Object



The University of Sydney

Amezquita et al. Nature Methods 2019   Page 25

We will try this soon …

14:15 – 15:00 Quality control and
data integration

THE UNIVERSITY OF
SYDNEY

# Roadmap for the workshop

Setting up: 13:30 – 13:45 Google cloud set up

Session 1: 13:45 – 14:15 Single cell analysis overview (scdney)

Session 2: 14:15 – 15:00 Quality control and data integration

AFTERNOON TEA: 1500-1530

Session 3: 15:30 – 16:00 Overview of single-cell downstream analysis

Session 4: 16:00 – 16:45 Downstream analysis: cell type identification, identify marker genes & cell type composition

Extension: cell type identification via supervised classification and single cell trajectory analysis

# Summary and Q&A

THE UNIVERSITY OF
SYDNEY

# Afternoon Tea

THE UNIVERSITY OF
SYDNEY

# Component 4: Cell type identification



Cell type identification

Unsupervised apparoach    Semi-supervised apparoach    Supervised apparoach

scReClassify

scClassify

## Science questions

- What cell types are present in the dataset?

- Can we identify the cell types?

# Component 4: Cell type identification



## Science questions

- What cell types are present in the dataset?

- Can we identify the cell types?

## Analysis techniques

- Visualization (dimension reduction)

- Clustering (unsupervised learning)

- Classification (supervised learning)

# Dimension reduced plot of our data (tSNE plot)

t–SNE plot



How many cell types are there?
What are the cell types?

# k-means clustering



t–SNE plot

# Clustering algorithms for scRNA-seq

*k*-means

Hierarchical

RaceID

SC3

CIDR

countClust

RCA

SIMLR



Luke Zappia, et al. *PLoS Comp. Bio.* 2018

# Similarity metric is the core of clustering algorithm

**Key question:** is there a similarity metric that performs (on average) better for clustering single cells based on their transcriptome?



k-means

Hierarchical

RaceID

SC3

CIDR

countClust

RCA

SIMLR

**Euclidean**

$$s_{ij} = \sqrt{\sum_{g=1}^{G} (x_{ig} - x_{jg})^2};$$

**Manhattan**

$$s_{ij} = \sum_{g=1}^{G} |x_{ig} - x_{jg}|;$$

**Maximum**

$$s_{ij} = \max_g |x_{ig} - x_{jg}|.$$

Distance-based

**Pearson**

$$s_{ij} = \frac{\sum_{g=1}^{G} (x_{ig} - \bar{x}_i)(x_{jg} - \bar{x}_j)}{\sqrt{\sum_{g=1}^{G} (x_{ig} - \bar{x}_i)^2} \sqrt{\sum_{g=1}^{G} (x_{jg} - \bar{x}_j)^2}};$$

**Spearman**

$$s_{ij} = \frac{\sum_{g=1}^{G} (r_{ig} - \bar{r}_i)(r_{jg} - \bar{r}_j)}{\sqrt{\sum_{g=1}^{G} (r_{ig} - \bar{r}_i)^2} \sqrt{\sum_{g=1}^{G} (r_{jg} - \bar{r}_j)^2}},$$

Correlation-based

# *k*-means Clustering on GSE60361

k-means



(a)

Annotated cells (GSE60361)

pre-defined cell types

- pyramidal CA1
- pyramidal SS
- interneurons
- microglia
- oligodendrocytes
- endothelial mural
- astrocytes ependymal

Zeisel A, et al. *Science* 2015

# Evaluation framework



Agreement to pre-defined classes:
Normalized Mutual Information (NMI)
Adjusted Rand Index (ARI)
Fowlkes-Mallows Index (FM)
Jaccard Index (Jaccard)

Taiyun Kim

# Evaluation results (against the pre-defined cell types)



## Impact of similarity metrics on single-cell RNA-seq data clustering

Taiyun Kim, Irene Rui Chen, Yingxin Lin, Andy Yi-Yang Wang,
Jean Yee Hwa Yang, Pengyi Yang

*Briefings in Bioinformatics*, bby076,

**PhD student: Taiyun Kim**

# Evaluation results (against the pre-defined cell types) using other measures



On average, correlation-based metrics improved on distance-based metrics by 31.5% (NMI), 39.6% (ARI), 16% (FM), 23% (Jaccard)

# Account for data scaling and zero-counts



| Source | Publication | Organism | # cell | # class |
|---|---|---|---|---|
| GSE45719 | Deng *et al.* (2014) | Mouse | 300 | 8 |
| GSE63818 | Guo *et al.* (2015) | Human | 328 | 37 |
| GSE67835 | Darmanis *et al.* (2015) | Human | 420 | 8 |
| GSE82187 | Gokce *et al.* (2016) | Mouse | 705 | 10 |
| GSE75140 | Camp *et al.* (2015) | Human | 734 | 13 |
| GSE75748 (TC) | Chu *et al.* (2016) | Human | 758 | 6 |
| GSE84133 | Baron *et al.* (2016) | Mouse | 822 | 13 |
| GSE89232 | Breton *et al.* (2016) | Human | 957 | 4 |
| GSE75748 (CT) | Chu *et al.* (2016) | Human | 1018 | 7 |
| GSE94820 | Villani *et al.* (2017) | Human | 1140 | 5 |
| E-MTAB-4079 | Scialdone *et al.* (2016) | Mouse | 1205 | 4 |
| GSE84371 | Habib *et al.* (2016) | Mouse | 1402 | 8 |
| GSE59114 | Kowalczyk *et al.* (2015) | Mouse | 1428 | 6 |
| E-MTAB-3929 | Petropoulos *et al.* (2016) | Human | 1529 | 5 |
| GSE93593 | Close *et al.* (2017) | Human | 1733 | 4 |
| GSE86146 | Li *et al.* (2017b) | Human | 2621 | 45 |
| GSE60361 | Zeisel *et al.* (2015) | Mouse | 3005 | 7 |
| GSE70630 | Tirosh *et al.* (2016b) | Human | 4347 | 8 |
| GSE72056 | Tirosh *et al.* (2016a) | Human | 4645 | 7 |
| Broad Portal | Habib *et al.* (2017) | Mouse | 13313 | 26 |
| Broad Portal | Habib *et al.* (2017) | Human | 14963 | 19 |
| GSE81905 | Shekhar *et al.* (2016) | Mouse | 27499 | 19 |

Data preparation
- Pre-processing
- Stratified subsampling

$X_1$
$X_2$
$\vdots$
$X_d$

Clustering algorithm
- Similarity metrics

Additional processing
- Linnorm normalisation
- SAVER imputation

Evaluation measures
- NMI
- ARI
- FM
- Jaccard

Pre-defined $P$ classes

$K$-means clustering with all or subset of genes
($K=P$)

**Agreement to pre-defined classes:**
Normalized Mutual Information (NMI)
Adjusted Rand Index (ARI)
Fowlkes-Mallows Index (FM)
Jaccard Index (Jaccard)

# Account for normalisation and imputation

# Improving the state-of-the-art clustering method using correlation metric

SIMLR

$$K(x_i, x_j) = \frac{1}{\epsilon_{ij}\sqrt{2\pi}} \exp\left(-\frac{}{2\epsilon_{ij}^2}\right)$$

$$s_{ij} = \frac{\sum_{g=1}^{G}(x_{ig} - \overline{x}_i)(x_{jg} - \overline{x}_j)}{\sqrt{\sum_{g=1}^{G}(x_{ig} - \overline{x}_i)^2}\sqrt{\sum_{g=1}^{G}(x_{jg} - \overline{x}_j)^2}};$$
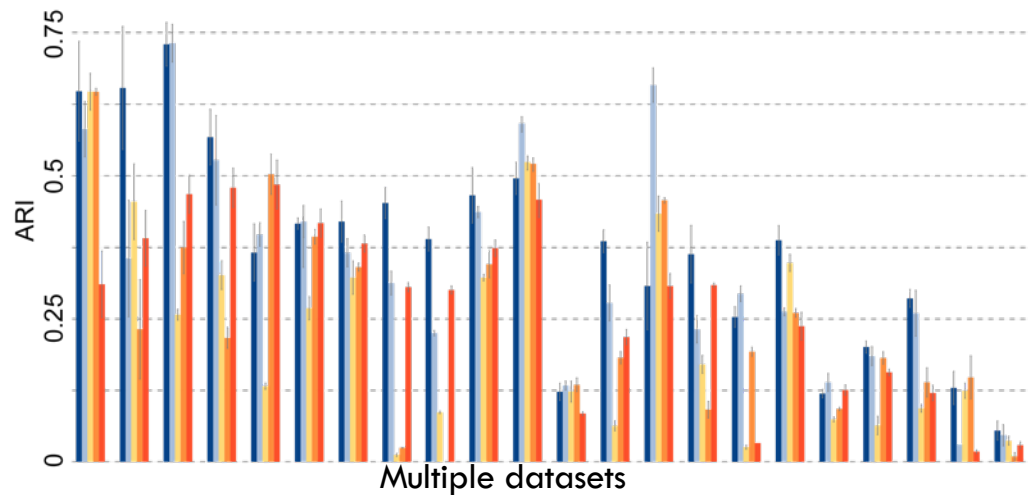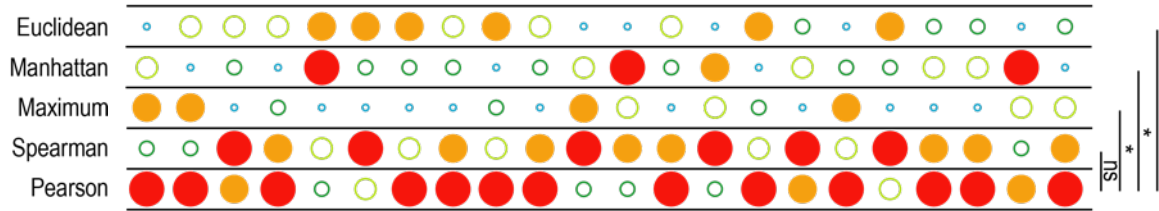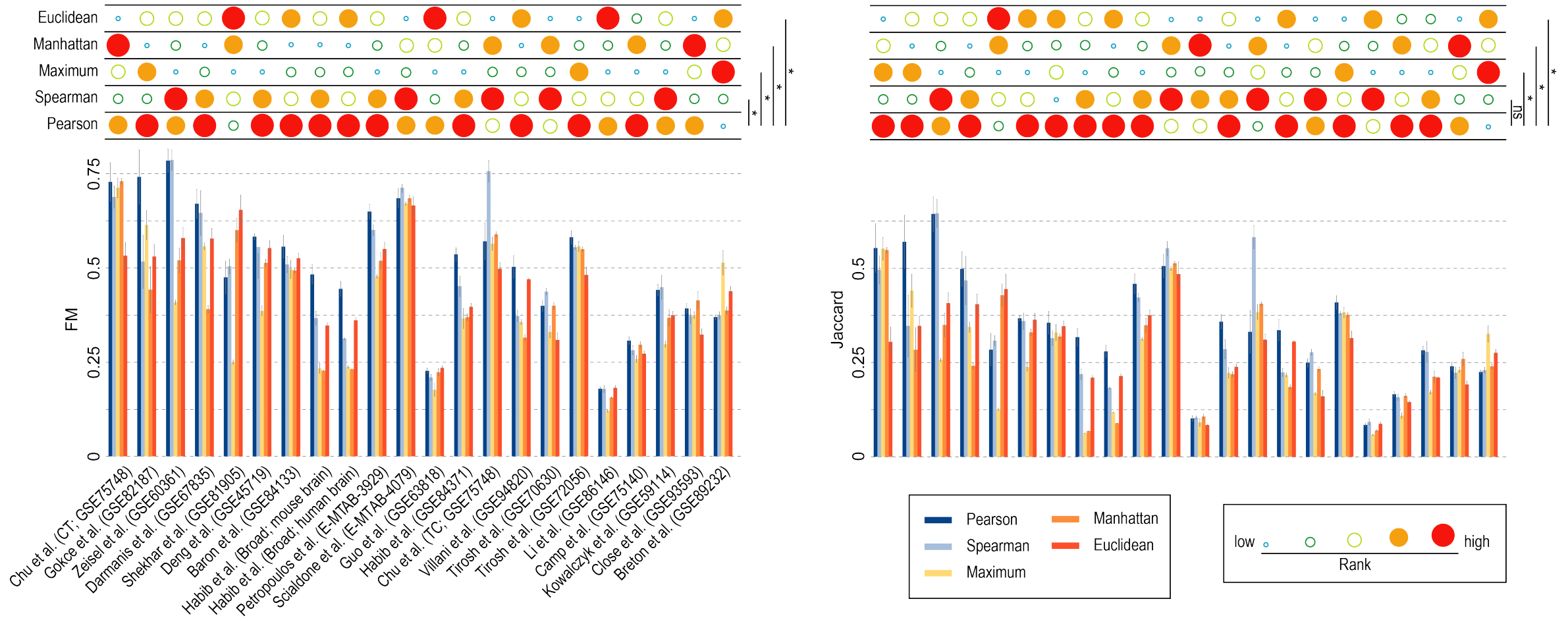


Wang, B., Zhu, J., Pierson, E., Ramazzotti, D., and Batzoglou, S. (2017). Visualization and analysis of single-cell rna-seq data by kernel-based similarity learning. *Nature Methods*, **14**(4), 414.

# Evaluation results of SIMLR with Pearson or Euclidean metrics

# Component 5: Downstream analysis



## Science questions

- Which genes are differentially expressed between cell types?

- What are the marker genes for each cell type?

- What is the cell type composition?

- Are the cells transitioning from one state to another?

# Differential expression testing: Differences between single cell and bulk RNAseq

- Advantage of single-cell:

    Account for cellular heterogeneity: DE tests can be now performed within cell-identity clusters across experimental conditions.

- Unique challenges for single-cell:

    - Dropout
    - High cell-to-cell variability

- Bulk DE methods

    - edgeR
    - limma
    - DESeq2

- Single-cell DE methods

    - MAST
    - ZINB-WaVE
    - DECENT
    - …

# DE methods comparisons for scRNAseq

# Cell type composition

Can we conclude that there are more cholangiocytes than mesenchymal cells?

# Single cell Differential Composition (scDC)

scDC simulates *uncertainty* in cell-type proportions via bootstrapping

Main components:
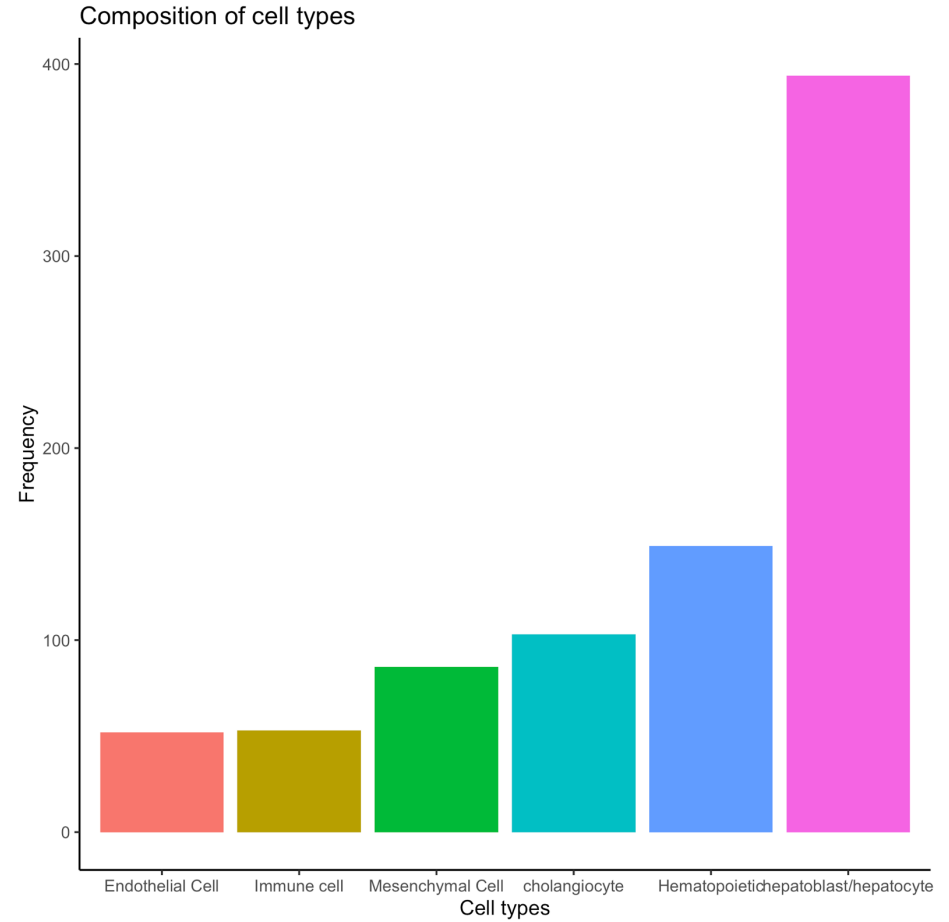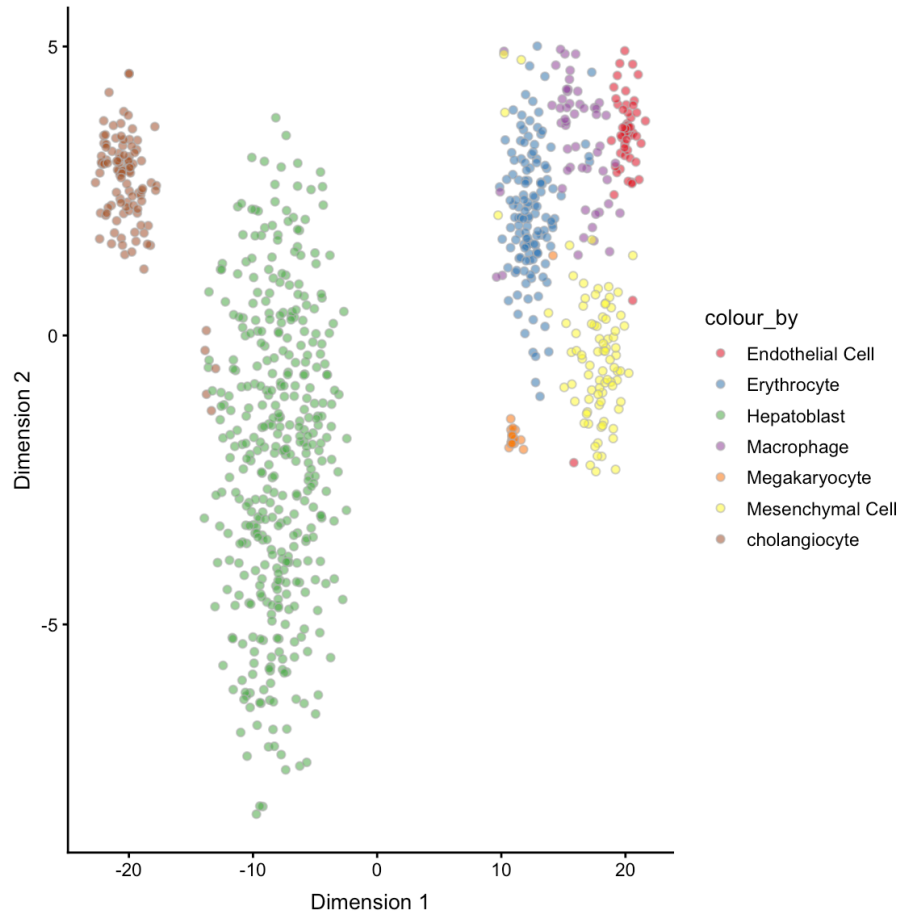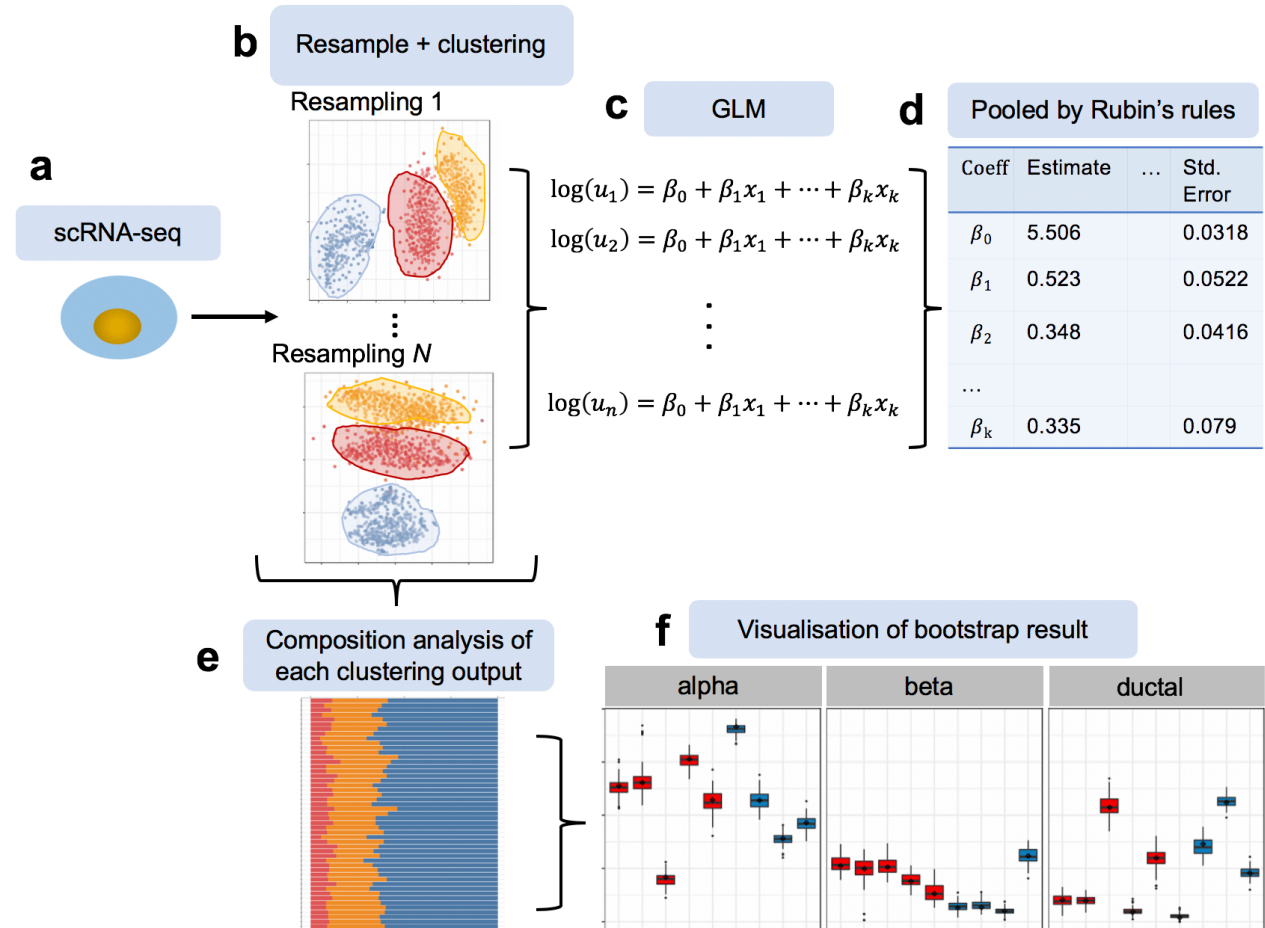- Sample with replacement from count matrix, stratified by patient
- Cell type identification via clustering (PCA -> Kmeans (Pearson correlation)
- Calculations of cell – type proportions standard error from bootstrap samples
- Calculation of pooled log-linear model using Rubin's pooled estimate



**a** scRNA-seq

**b** Resample + clustering

Resampling 1

Resampling *N*

**c** GLM

$$\log(u_1) = \beta_0 + \beta_1 x_1 + \cdots + \beta_k x_k$$
$$\log(u_2) = \beta_0 + \beta_1 x_1 + \cdots + \beta_k x_k$$
$$\vdots$$
$$\log(u_n) = \beta_0 + \beta_1 x_1 + \cdots + \beta_k x_k$$

**d** Pooled by Rubin's rules

| Coeff | Estimate | … | Std. Error |
|---|---|---|---|
| $\beta_0$ | 5.506 | | 0.0318 |
| $\beta_1$ | 0.523 | | 0.0522 |
| $\beta_2$ | 0.348 | | 0.0416 |
| … | | | |
| $\beta_k$ | 0.335 | | 0.079 |

**e** Composition analysis of each clustering output

**f** Visualisation of bootstrap result

| alpha | beta | ductal |

# Single cell Differential Composition (scDC)

– Examined two synthetic datasets constructed from two sets of real experimental data — Pancreas (T2D vs healthy) and Neuronal (developing mouse)

– In pancreas dataset

- confirmed the original finding that 1 of the 4 subjects has a higher beta cell value, as IQR non overlap

– In neuronal dataset

- Revealed new finding that progenitor cells percentage increase over time



**d** Cell Proportion Bootstrap Distribution

**e** Average Cell Proportion Across Subjects

**f** Mean Cell Composition by Bootstrap

We will try this soon…

16:00 – 16:45 Downstream analysis: identify marker genes & cell type composition
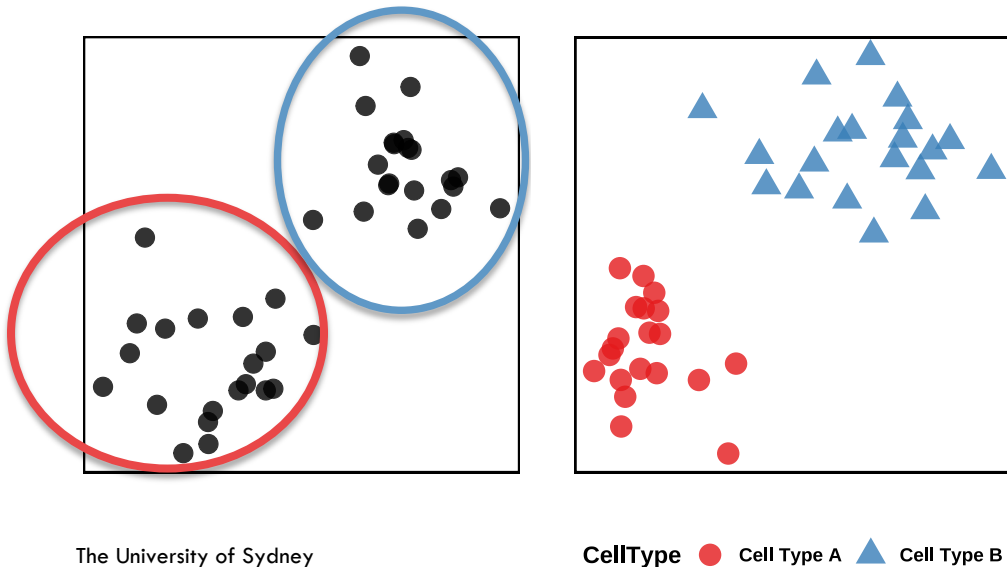
THE UNIVERSITY OF
SYDNEY

**Extension:**
**1. cell type identification via supervised classification**
**2. single cell trajectory analysis**

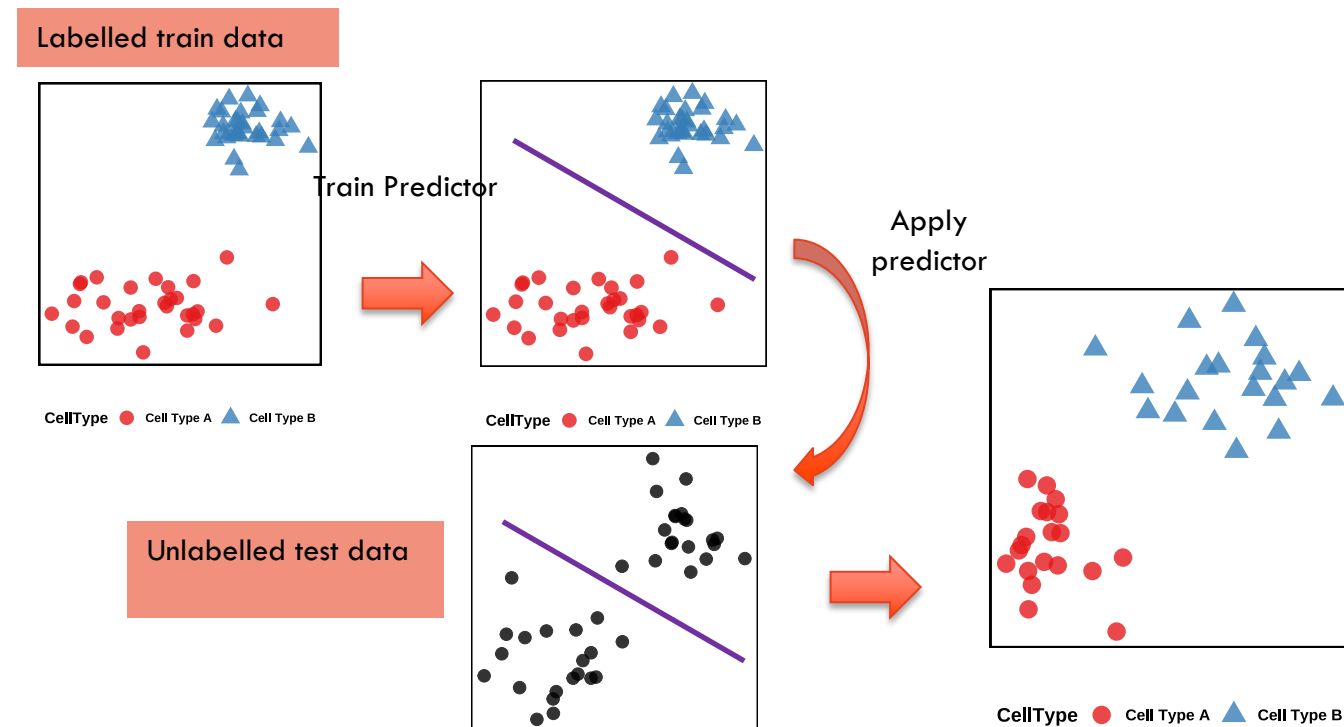# An alternative approach of cell type identification: supervised learning

## Clustering (unsupervised learning)

- Group the cells that are "close" to each other
- Annotated each cluster by DE genes or other characteristics
- Identify the novel cell type



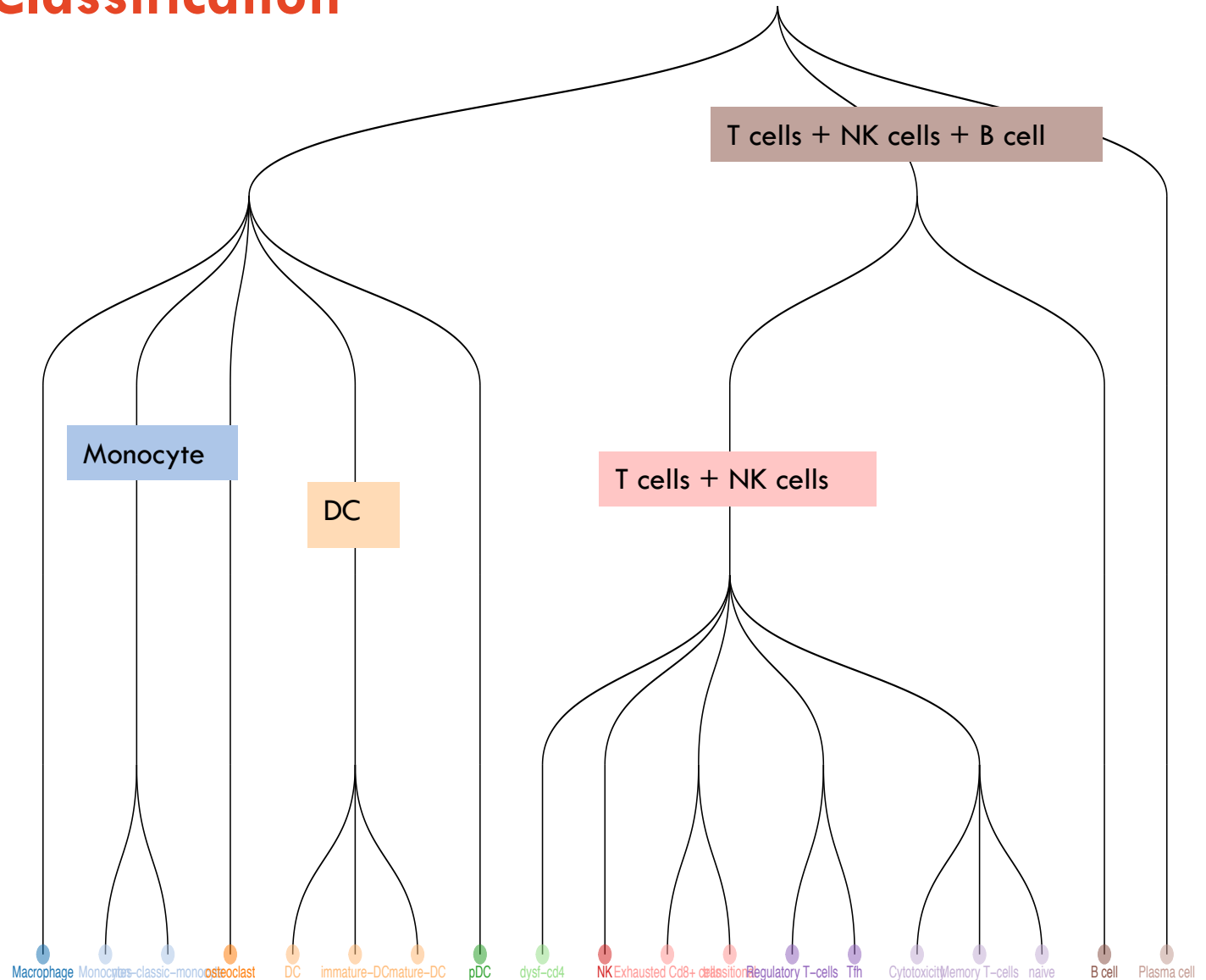CellType ● Cell Type A ▲ Cell Type B

## Classification (supervised learning)

- Required reference labelled datasets
- Predict cell types label directly
- What if there are cell types that are not in the reference data?

Labelled train data

Train Predictor

Apply predictor



CellType ● Cell Type A ▲ Cell Type B

CellType ● Cell Type A ▲ Cell Type B

Unlabelled test data

CellType ● Cell Type A ▲ Cell Type B

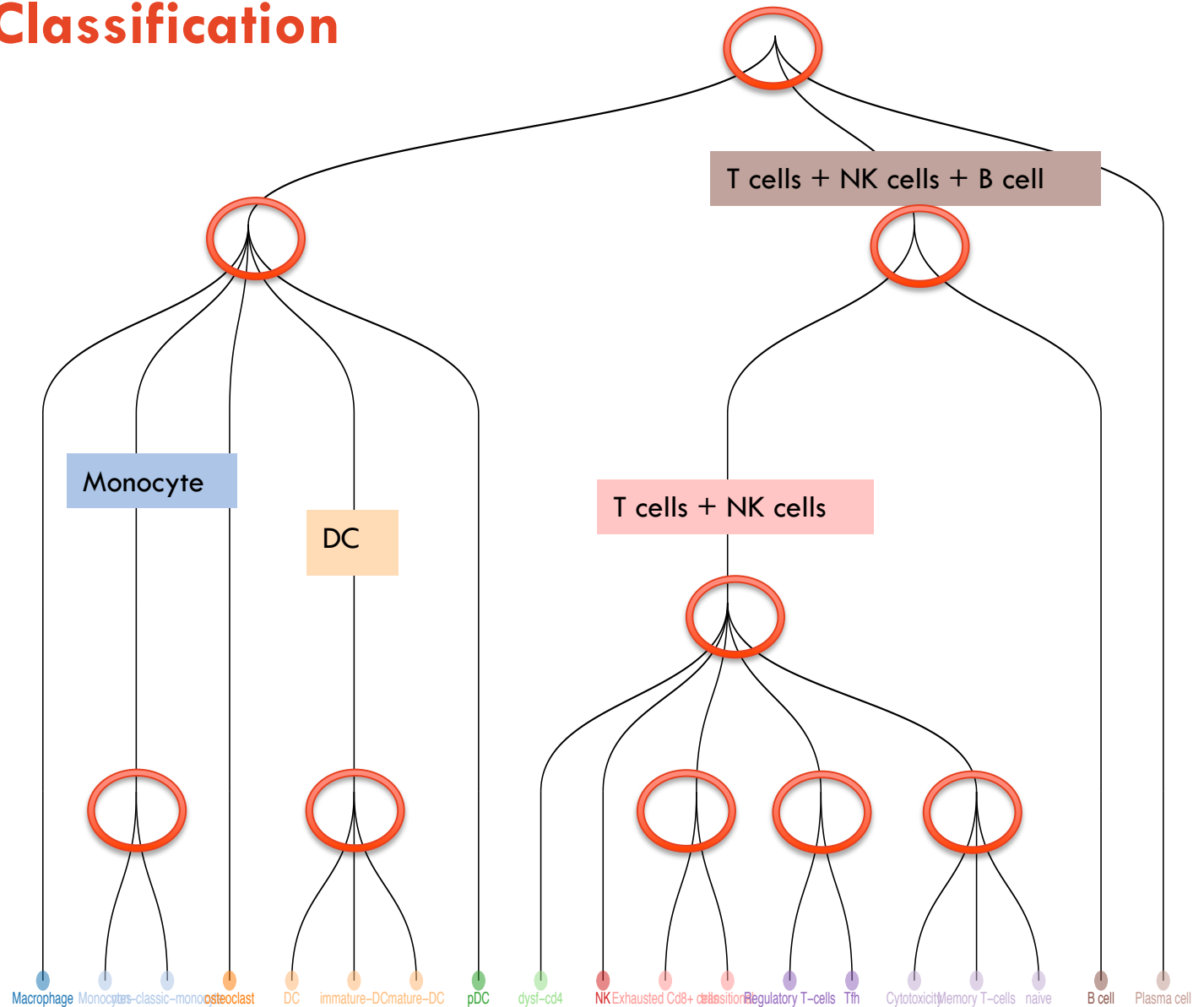# scClassify: Hierarchical Classification

**Step 1: Constructing cell type hierarchical tree:**

*We use hierarchical ordered partitioning and collapsing hybrid (HOPACH) to generate the cell type hierarchical tree based on the reference dataset.*
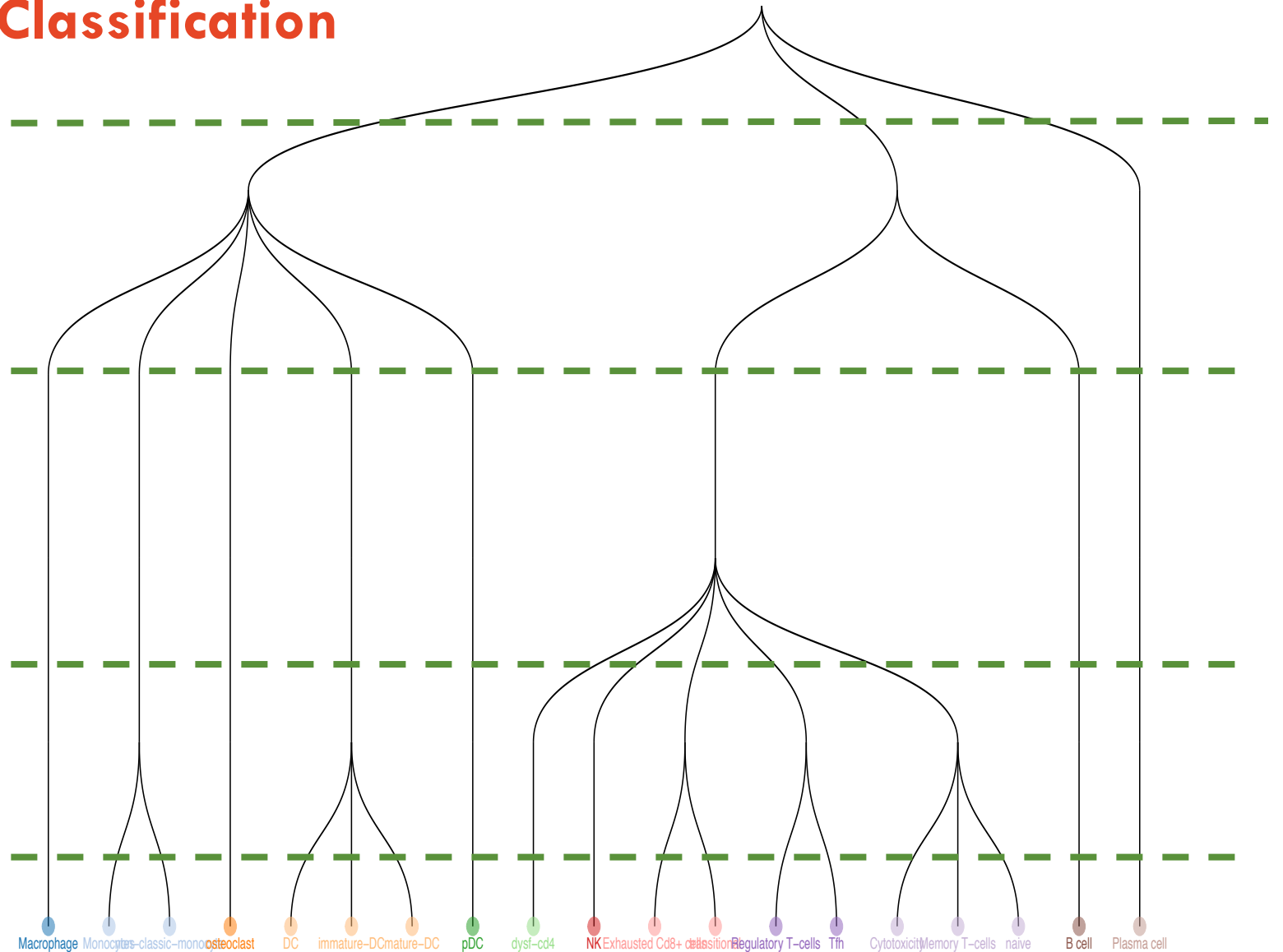
van der Laan, M. J. and Pollard, K. S. (2003), *Journal of Statistical Planning and Inference.*

# scClassify: Hierarchical Classification

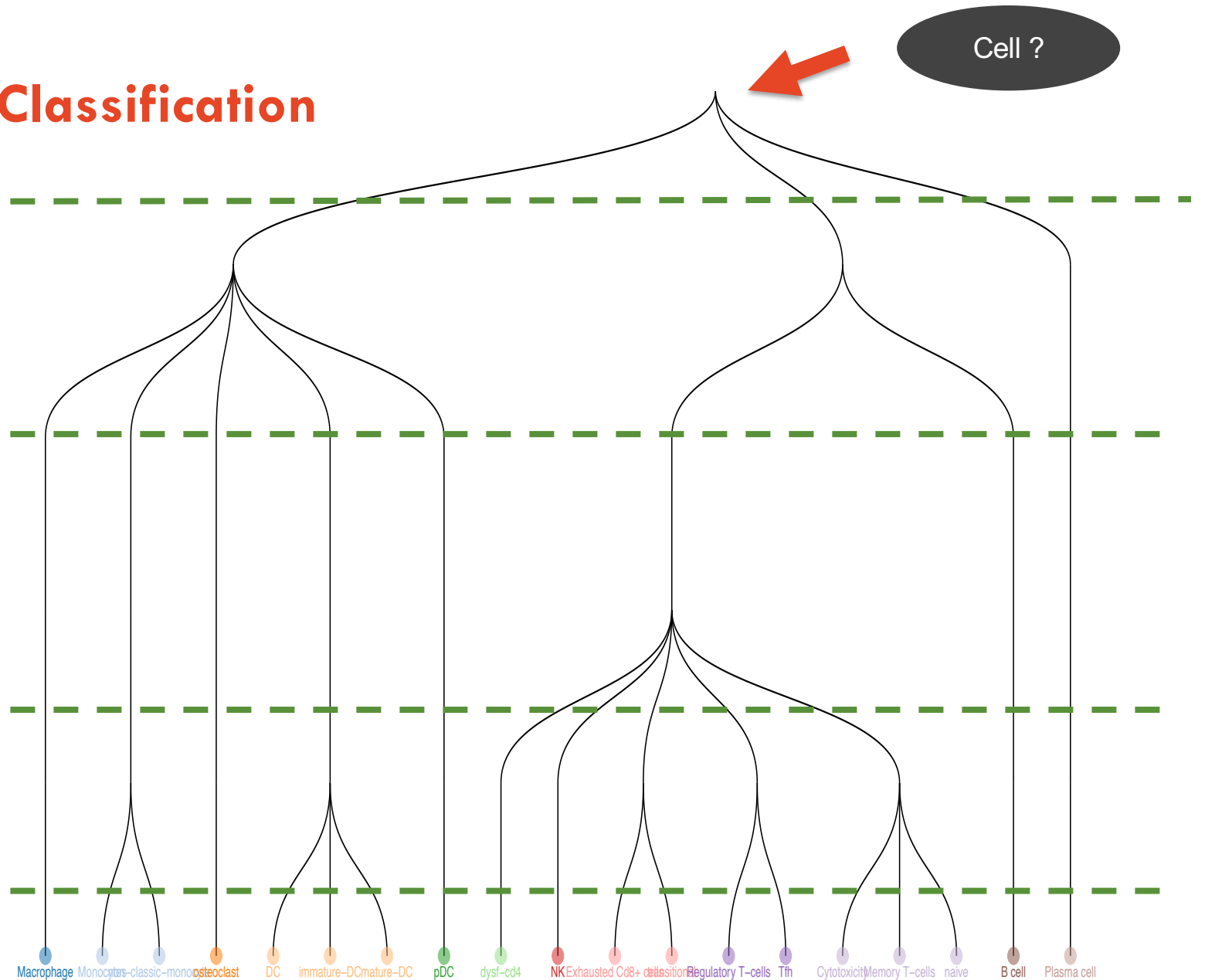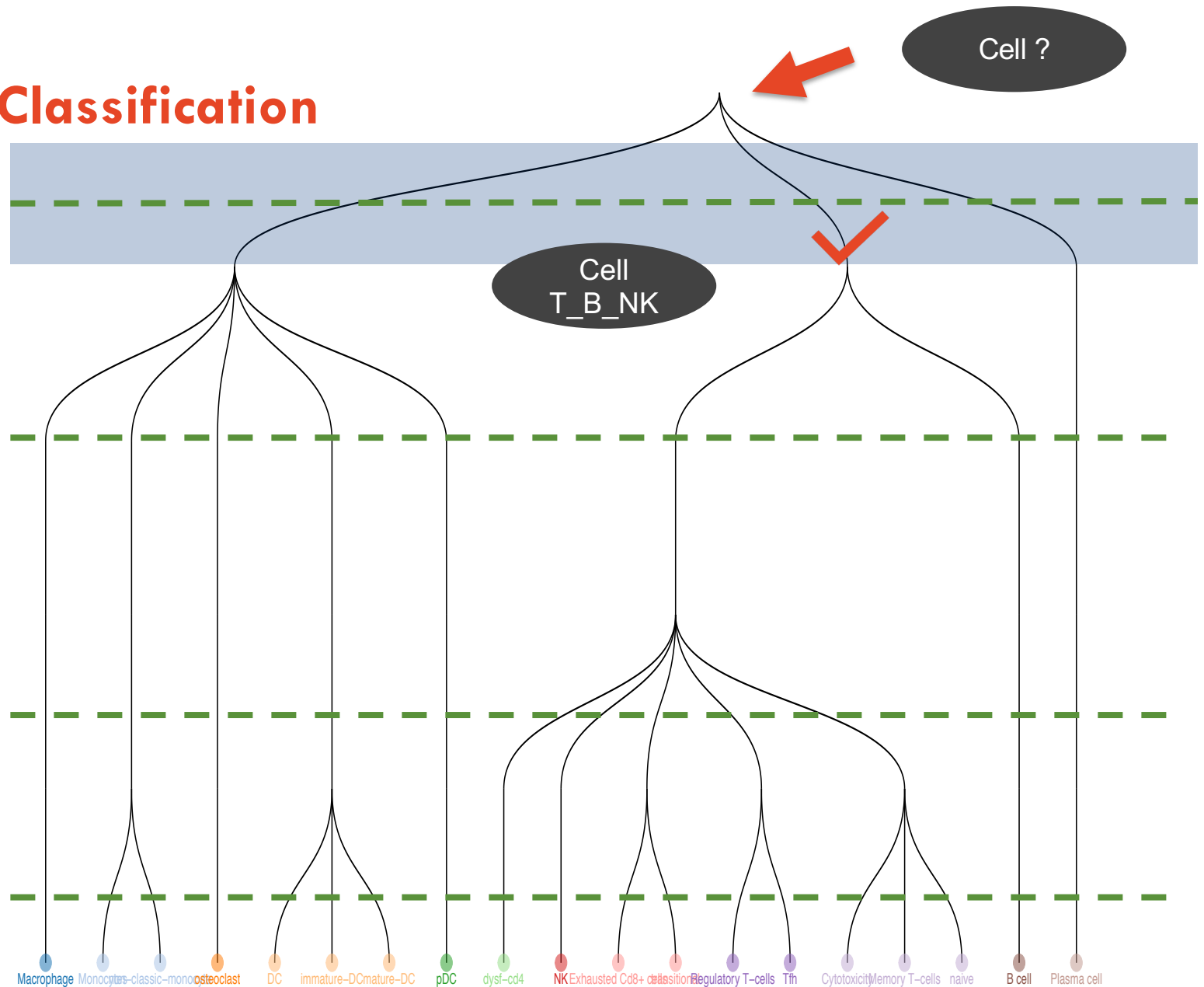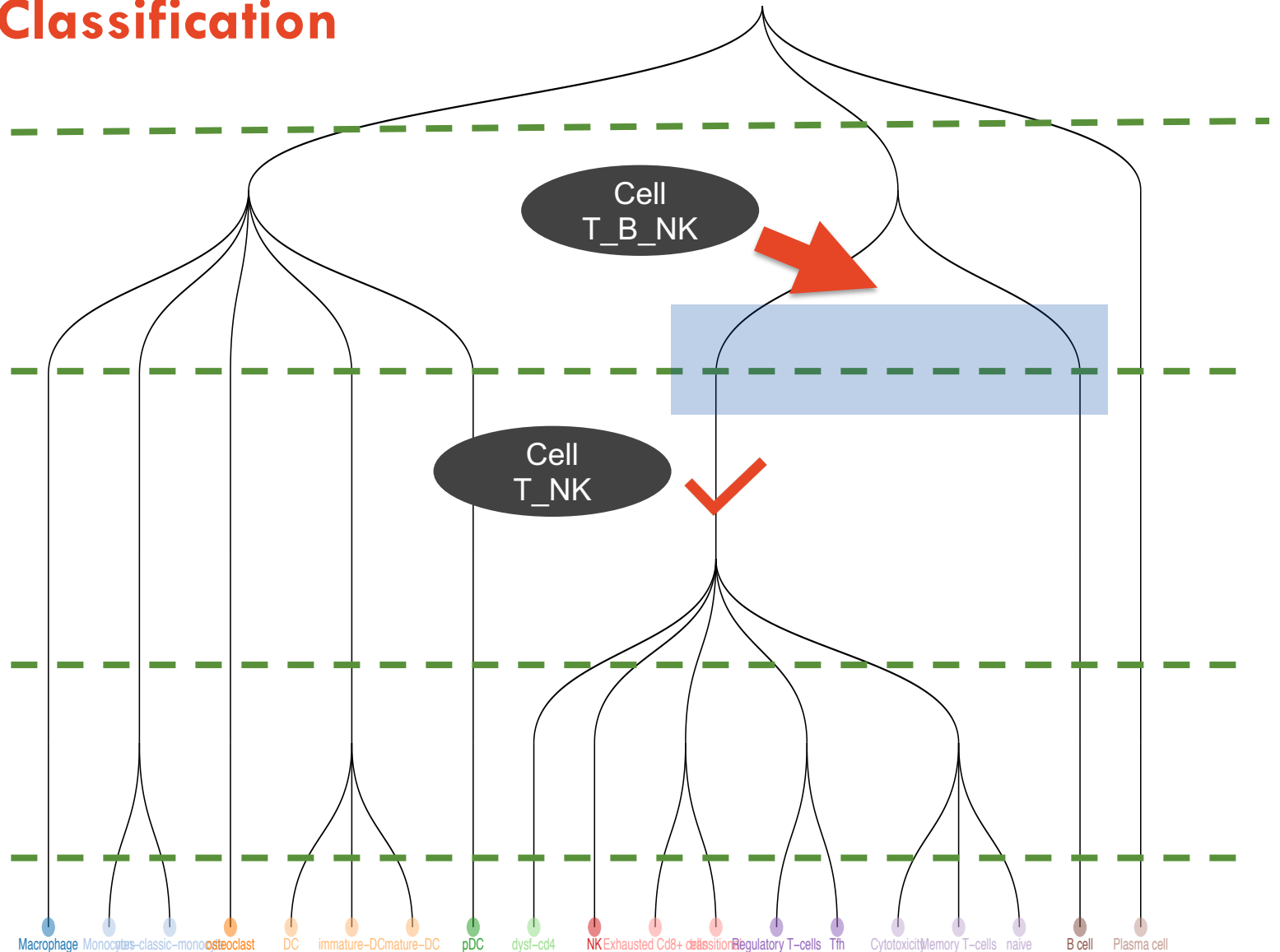*Step 2: Feature selection at each branch point.*

# scClassify: Hierarchical Classification

*Step 3: Performing correlation-based weighted kNN for each level of the cell type hierarchical tree:*

# scClassify: Hierarchical Classification

Cell ?

*Step 3: Performing correlation-based weighted kNN for each level of the cell type hierarchical tree:*

Macrophage  Monocytes-classic-monocyte  osteoclast  DC  immature-DC mature-DC  pDC  dysf-cd4  NK Exhausted Cd8+ cells transition Regulatory T-cells  Tfh  Cytotoxicity Memory T-cells  naive  B cell  Plasma cell
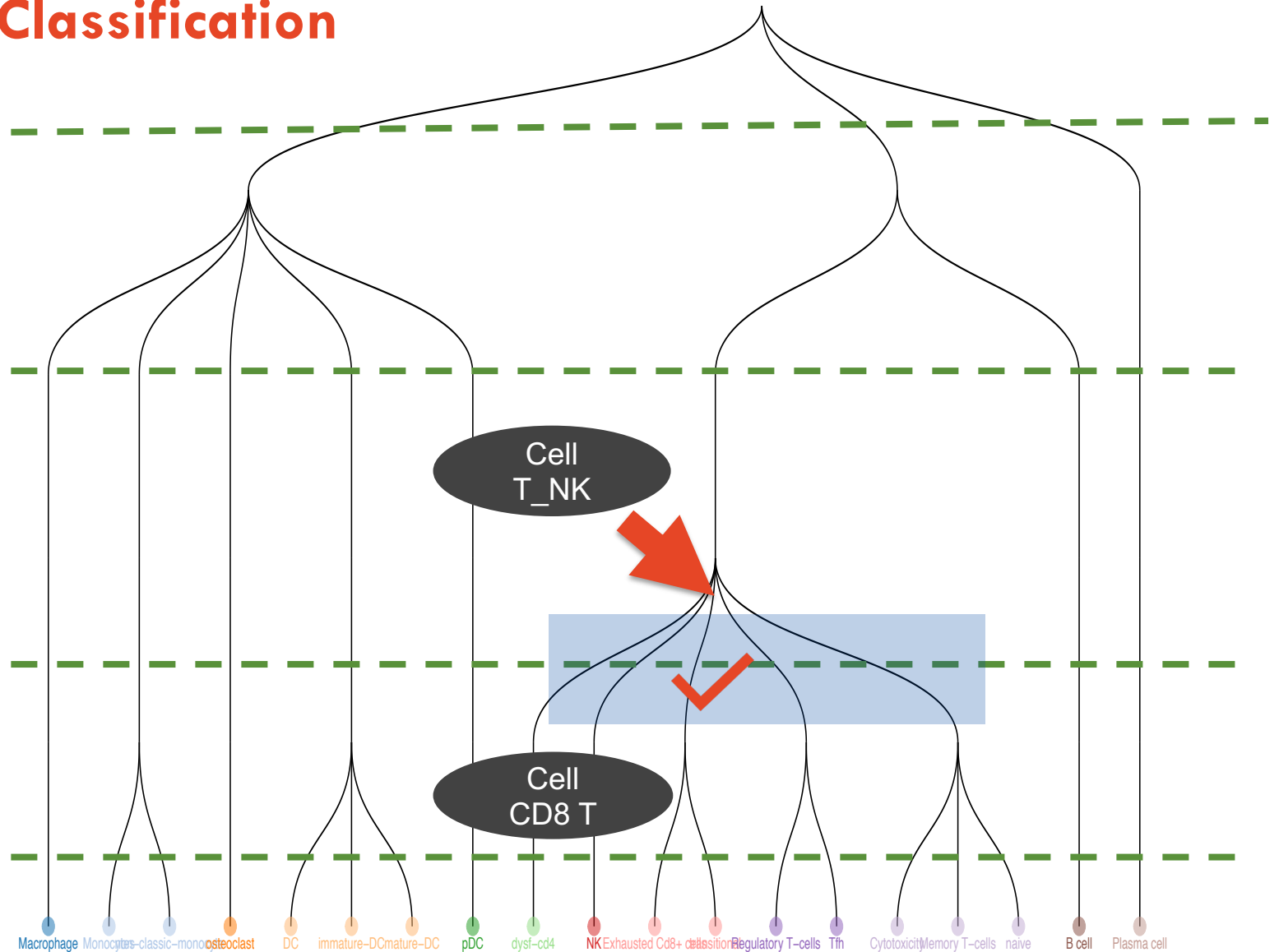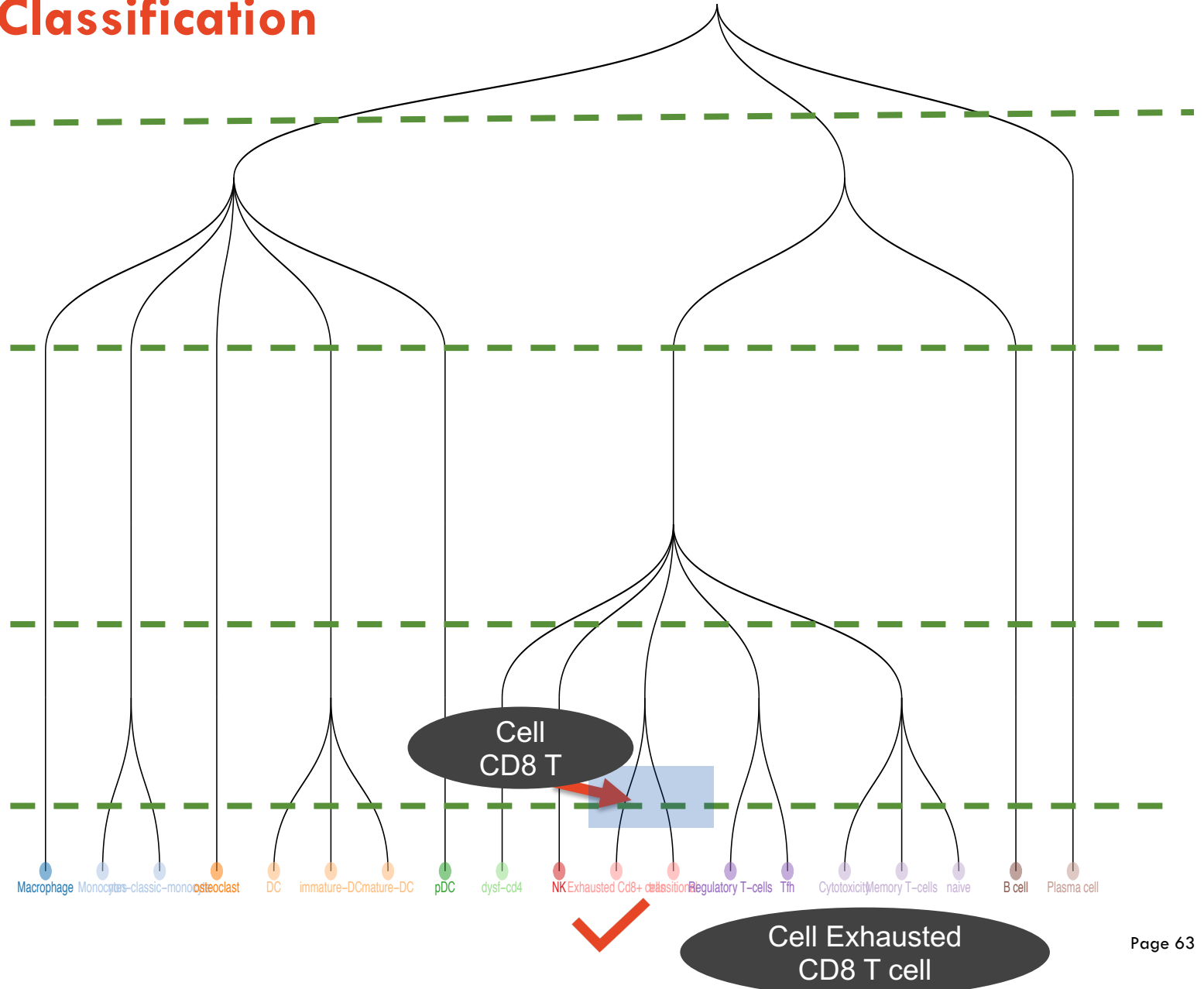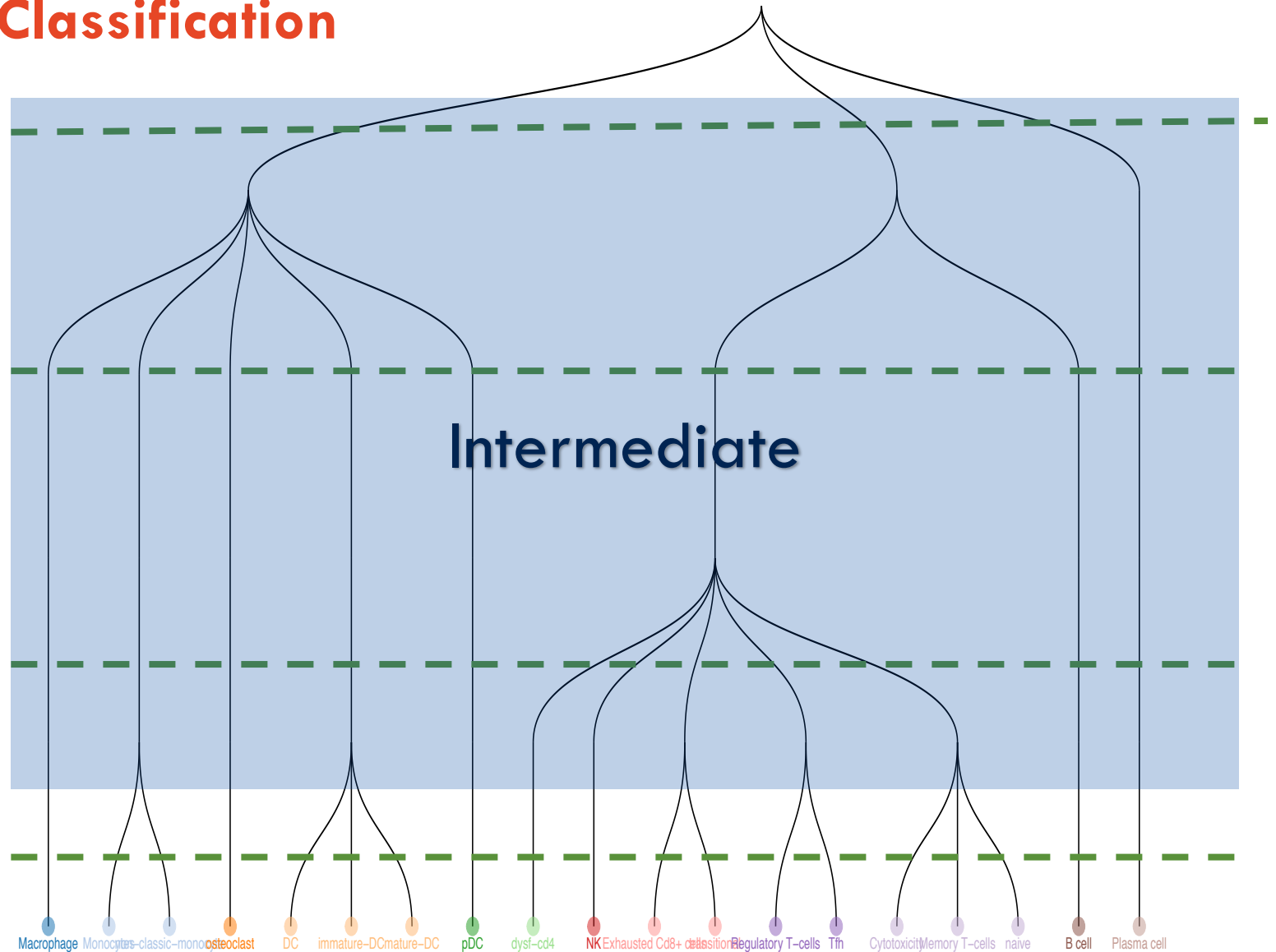
# scClassify: Hierarchical Classification

*Step 3: Performing correlation-based weighted kNN for each level of the cell type hierarchical tree:*

# scClassify: Hierarchical Classification

*Step 3: Performing correlation-based weighted kNN for each level of the cell type hierarchical tree:*

# scClassify: Hierarchical Classification

*Step 3: Performing correlation-based weighted kNN for each level of the cell type hierarchical tree:*
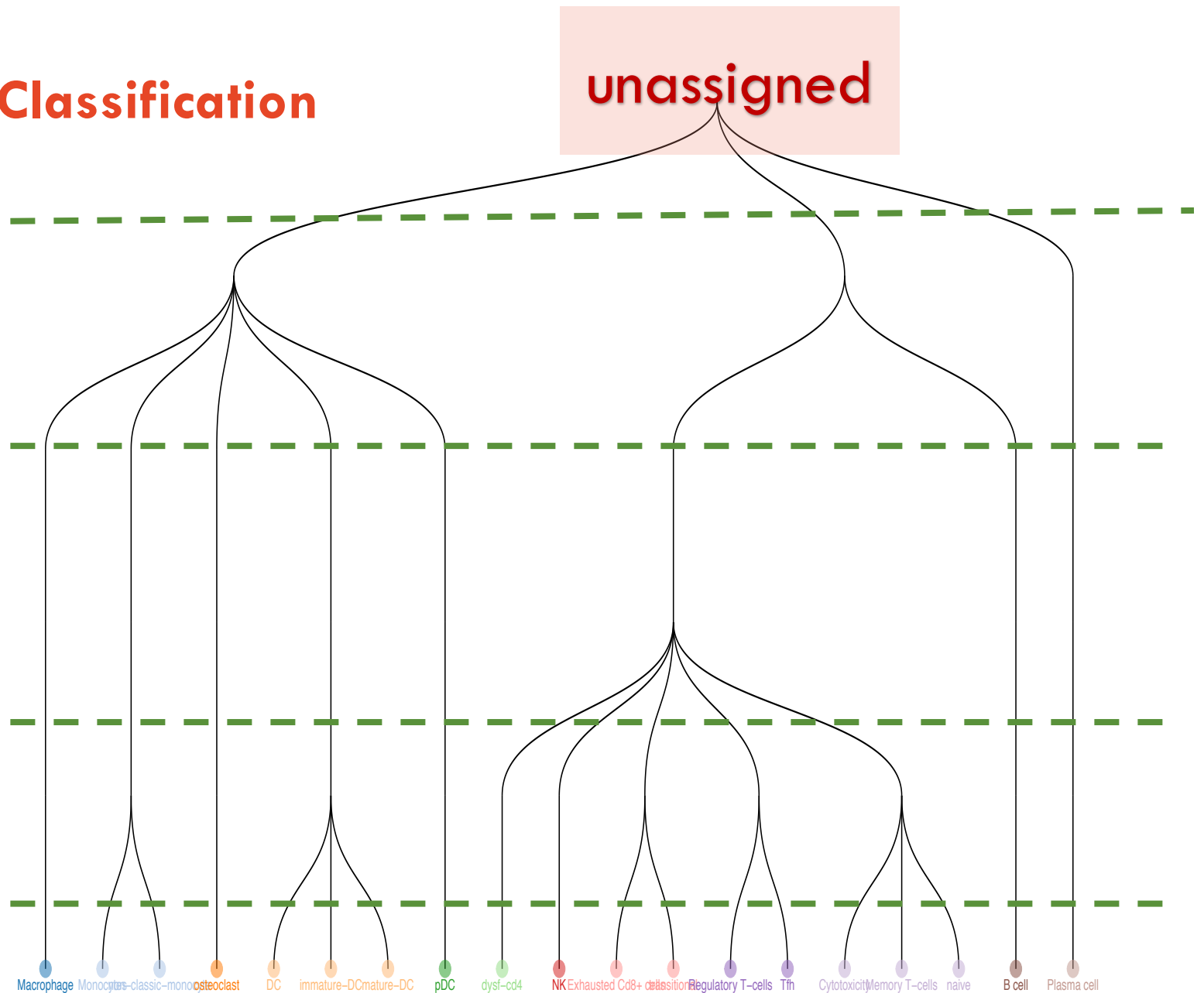
# scClassify: Hierarchical Classification

*Step 3: Performing correlation-based weighted kNN for each level of the cell type hierarchical tree:*

# scClassify: Hierarchical Classification

*Step 3: Performing correlation-based weighted kNN for each level of the cell type hierarchical tree:*

# scClassify: Hierarchical Classification

unassigned

*Step 3: Performing correlation-based weighted kNN for each level of the cell type hierarchical tree:*

Macrophage  Monocytes-classic-monocyte  osteoclast  DC  immature-DC  mature-DC  pDC  dysf-cd4  NK  Exhausted  Cd8+ cells  transition  Regulatory T-cells  Tfh  Cytotoxicity  Memory T-cells  naive  B cell  Plasma cell

# scClassify

Try scClassify: https://sydneybiox.github.io/scClassify/

# Trajectory inference

## Why trajectory analysis?

- Cells may not be sufficiently be described by a discrete classification system such as clustering
- Biological processes drive the development are usually continuous process
- Trajectory inference therefore can be used to model
  - the transitions between cell identities
  - Branching differentiation process
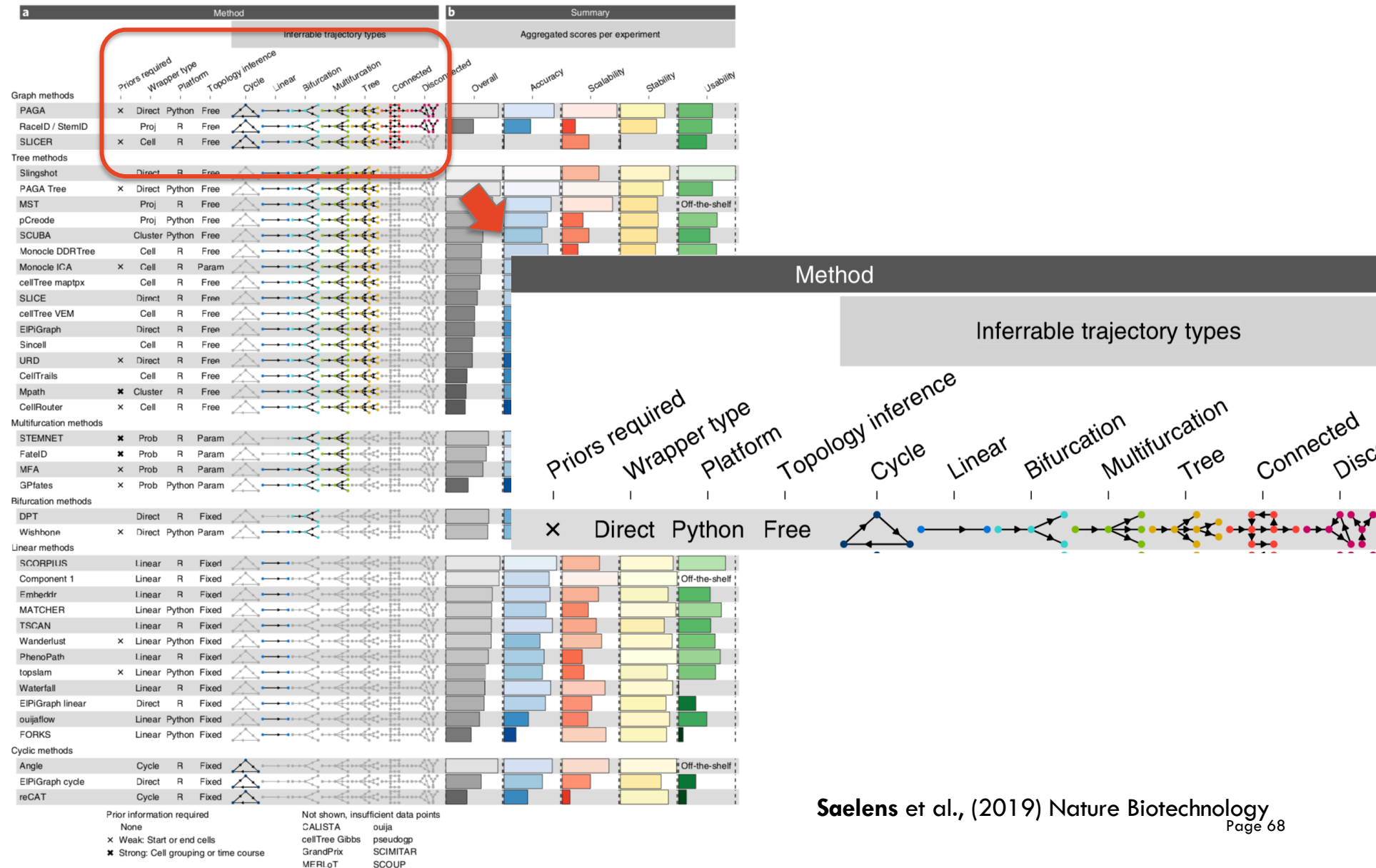  - Dynamic gene regularization model

## What is trajectory inference?

- Interpret single-cell data as a snapshot of a continuous process.

## Typical steps involved in trajectory inference:

- Reduce the dimensionality of the single cell data
- Finding paths through the reduced dimension space, by minimizing the changes between neighboring cells
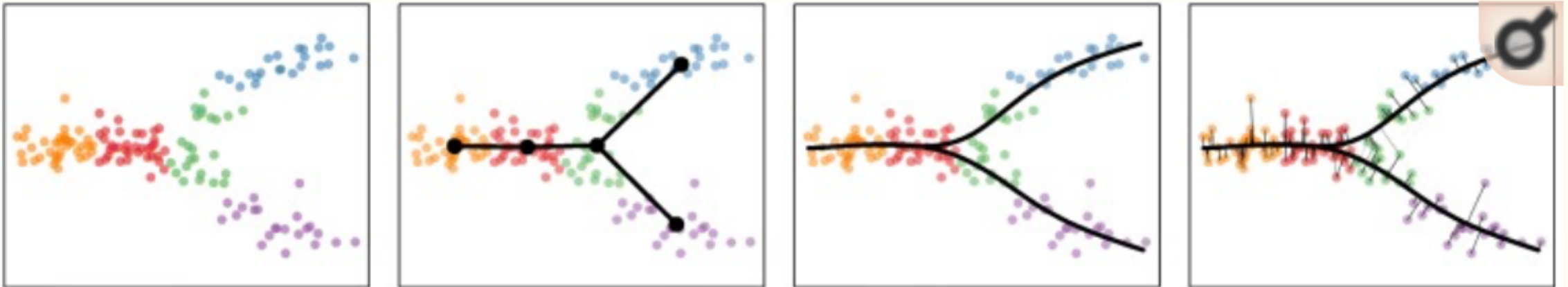- Order the cells by pseudotime

# Comparisons of pseudotime inference methods



**Saelens** et al., (2019) Nature Biotechnology

# Slingshot example (Street et al., 2018)

Three stages:

1. Reduced dimension of the data

2. Inference of the global lineage structure. Uses cluster-based minimum spanning tree

3. Inference of pseudotime variables for cells along each lineage. Fits simultaneous principal curves

# Summary

# Acknowledgment

### *The University of Sydney*

### *School of Mathematics and Statistics*

Sydney Precision Bioinformatics Research Group