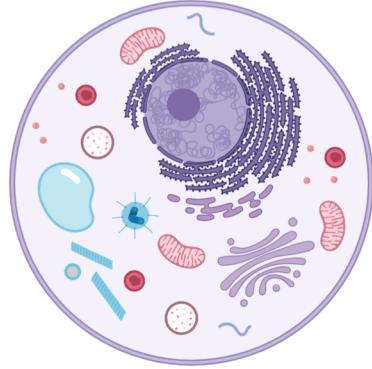


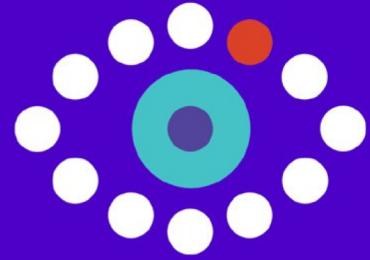


**Garvan Institute**  
of Medical Research



# Cell-type classification at single-cell resolution

Jose Alquicira Hernandez



# Outline

1. What is a cell type?
2. Difference between cell type and cell state
3. What is classification? What is prediction?
4. Types of classification
  - a. Uni vs. multivariate
  - b. Hierarchical vs. linear
  - c. Cluster-based vs. supervised
5. Software
6. Recommendations
7. scPred



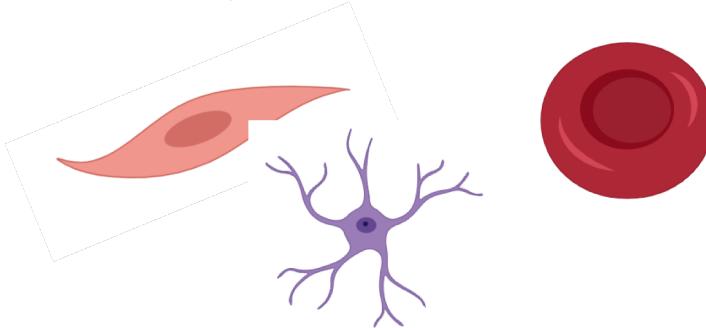
**Garvan Institute**  
of Medical Research

# What is a cell type?



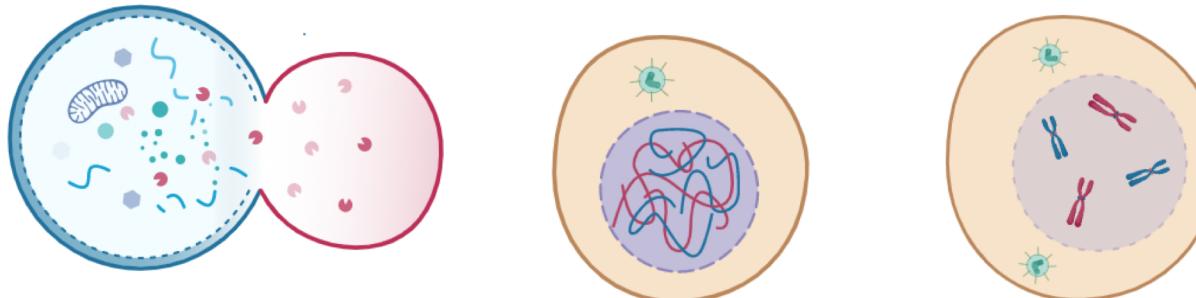
# Cell types can be defined according to different criteria

- Morphology
- Function
- Molecular composition
- End or start point of differentiation
- Determined fate
- *Evolutionary conservation*
- *Phenotype*



# Cell states are different to cell types

- Temporary instance of a cell type
- *Two cells can share a cell state and being distinct cell types*
- Cell states are more overlapped than cell types
- No determined fate



# Cell types can be detected using various methods



## Methods

- Transcriptome
  - Epigenome
  - Proteome
  - *Surface markers*
  - Metabolome
  - Morphology
  - *Imaging*
  - Spatial transcriptomics
- ...

## When are two cells the same?

- Depends on the question
- Context dependent
- Granularity
- Hierarchical definition using established cell ontology

## Cell type classification

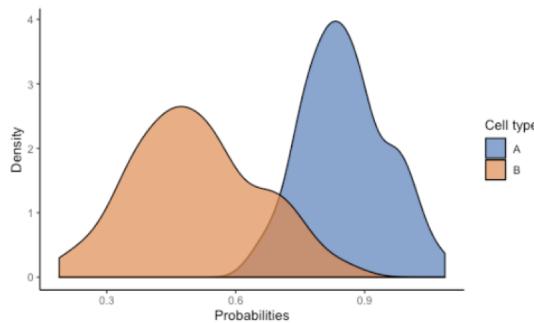
- Process to label a cell based on a previous definition and a set of features

- **Cell type definition**
- Ground truth based on biological and/or statistical criteria
- Cell type discovery

# Classification vs. prediction

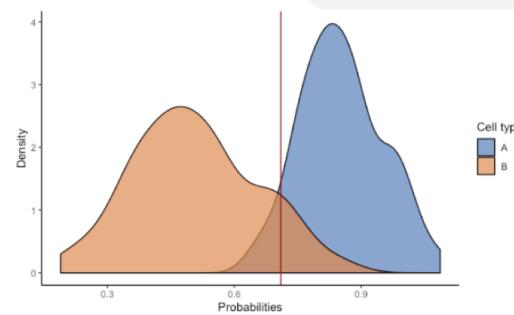
## Prediction

- Modeling of tendencies: probabilities

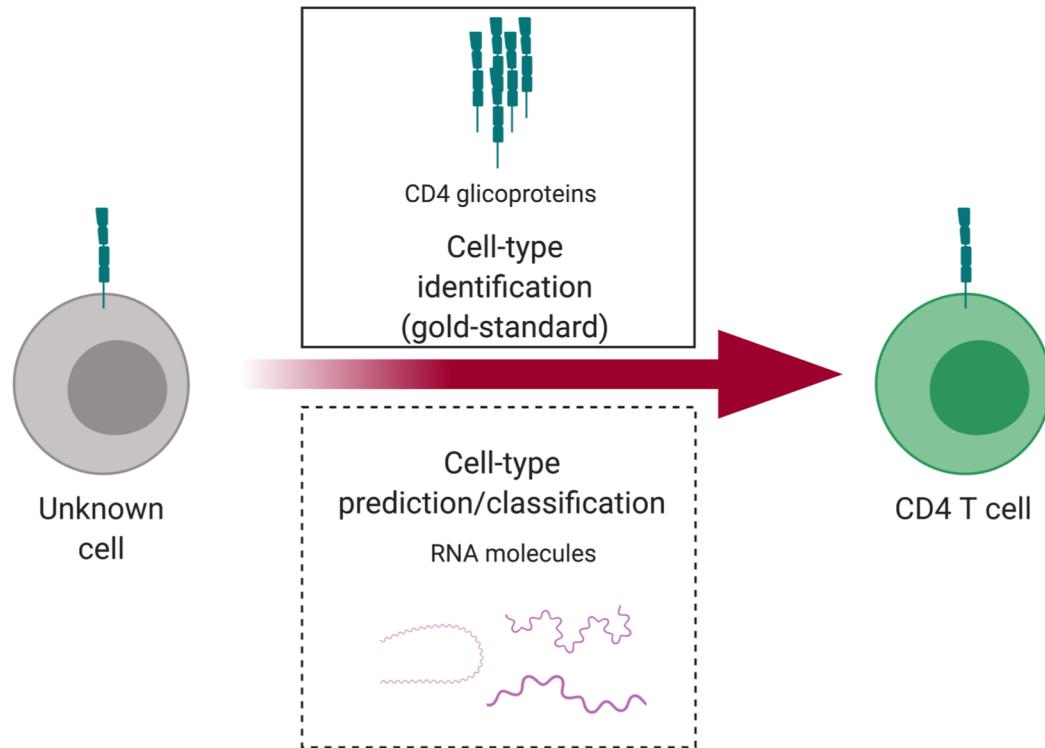


## Classification

- Decision making
- Forced choice according to context
- Classification combines prediction and decision making



Cell type classification is based on the premise that a set of features (e.g. gene expression) is able to recapitulate the variance of the phenotype we are interested in





**Garvan Institute**  
of Medical Research

# Cell-type classification approaches

# Univariate vs. multivariate

## Univariate

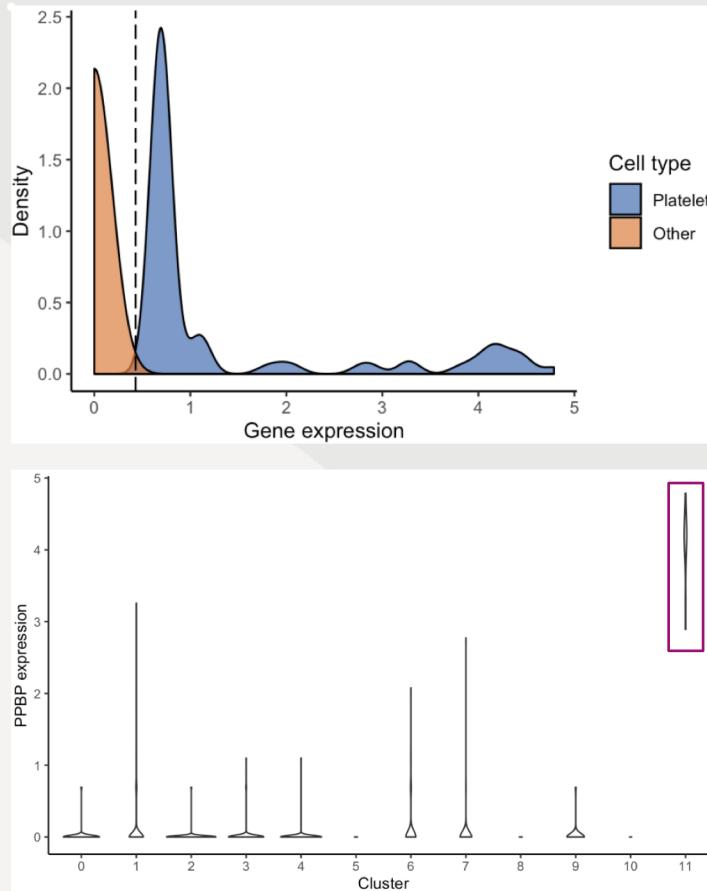
- Some cell types can be classified using a single gene marker (e.g. platelets) or protein

## Advantages

- Easier classification
- Clear interpretation

## Caveats

- Context dependent (shared expression between cells)
- Lack of correlation between canonical marker and feature
- Expression variance



# Univariate vs. multivariate

## Multivariate

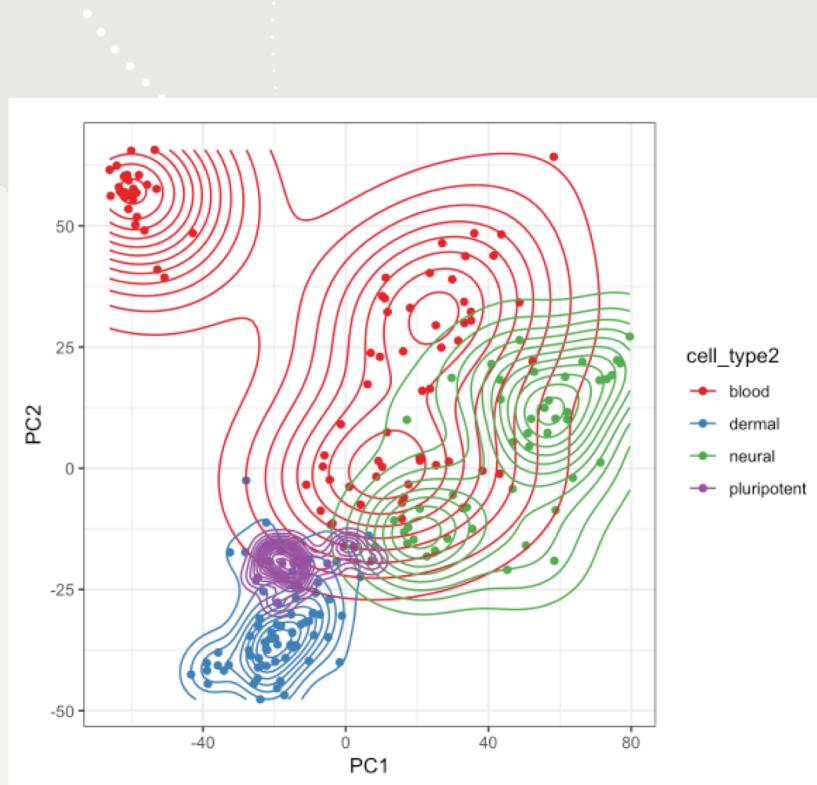
- Combination of features explain cell identity

## Advantages

- More information is used to classify a cell type (coexpression)

## Caveats

- Feature selection (DEGs, classic markers)



# Flat vs hierarchical

## Flat

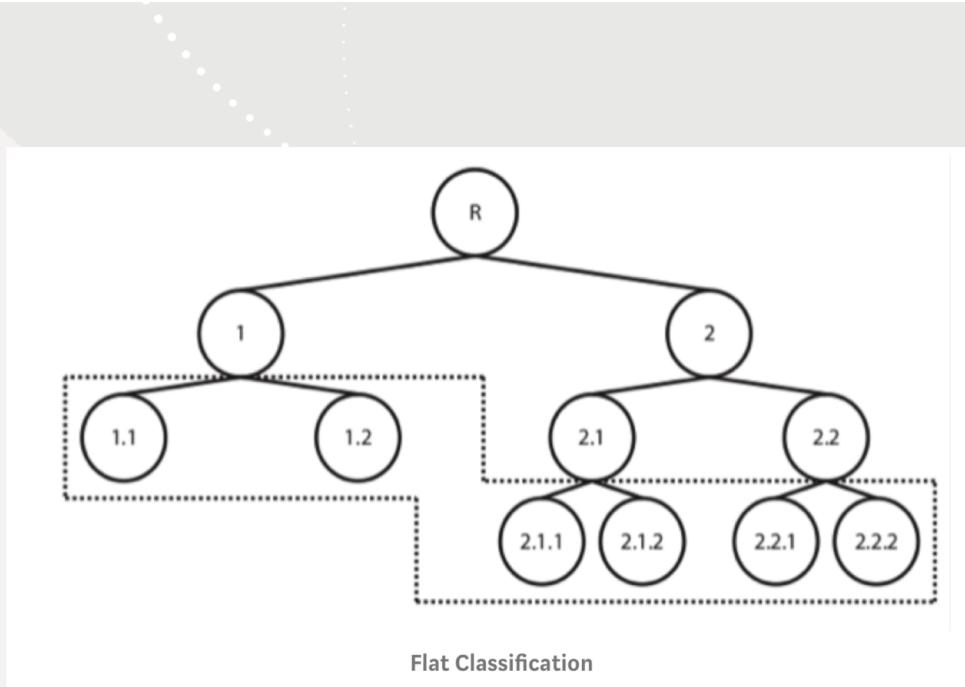
All cells are classified in a single step

## Advantages

- Simple
- Fast

## Caveats

- Cell heterogeneity (outlier populations)
- Cell type relatedness



From: "A Survey of Hierarchical Classification Across Different Application Domains"

# Linear vs hierarchical

## Hierarchical

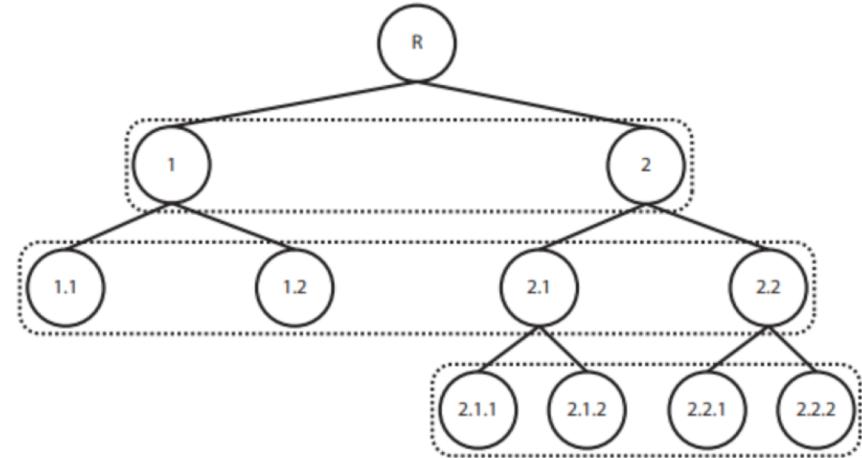
Takes into account cell organization/relatedness (e.g. hematopoietic lineage)

## Advantages

- Based on biological knowledge of the population

## Caveats

- Slower depending on the complexity of the hierarchy



Local Classifier Per Level Approach

From: "A Survey of Hierarchical Classification Across Different Application Domains"

# Unsupervised classification



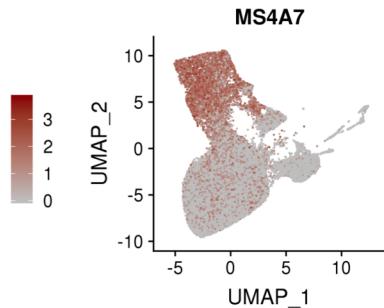
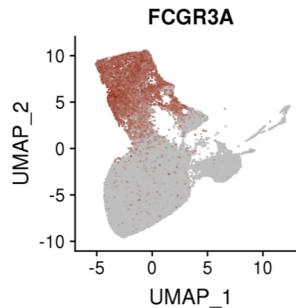
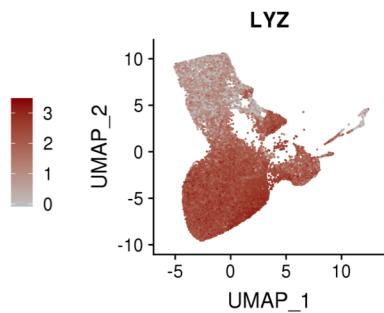
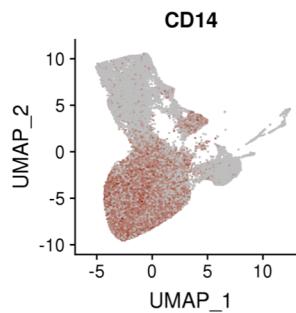
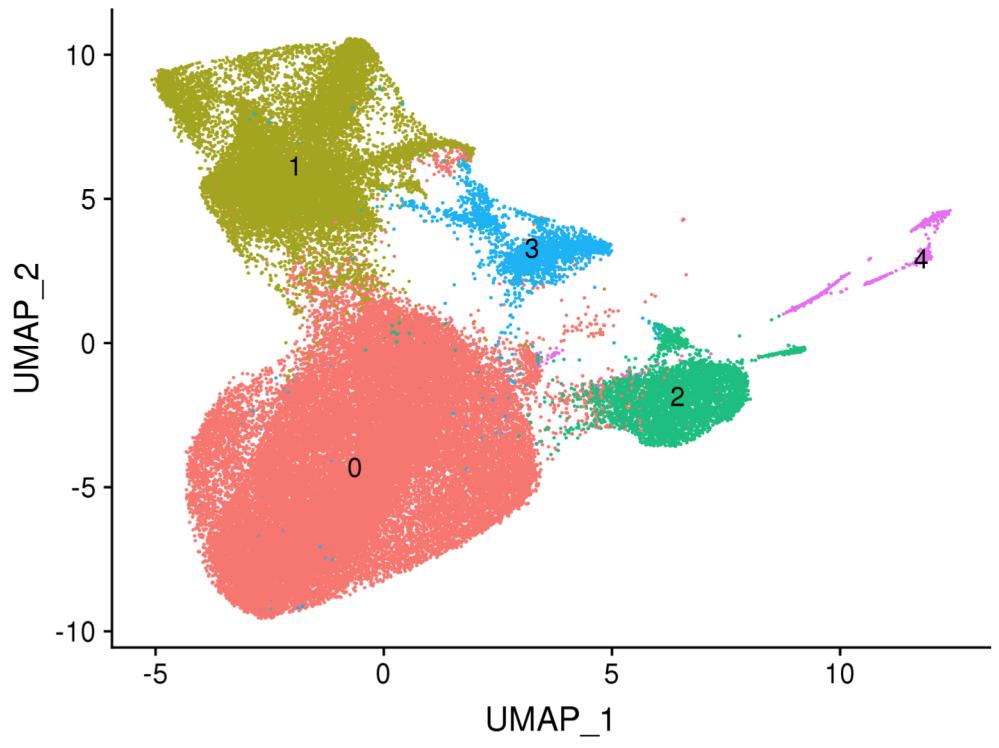
- Clustering
- Clusters are expected to correlate with real cell types or subpopulations
- All cells within a cluster are assumed to be the same cell type

## Advantages

- No reference is needed
- Classic markers can be used to guide the annotation
- Cell type discovery / definition

## Caveats

- Ambiguous results depending on resolution or hierarchy level
- Feature selection



# Supervised classification



A training dataset is used as reference to guide the classification of cells in the population of interest

## Advantages

- Fast to apply once the reference is built
- Classification performance estimated in training step
- Consistent classification criteria applied to different datasets

## Caveats

- Lack of reference (gold standard data)
- Completeness of reference

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

Tamim Abdelaal<sup>1,2#</sup> ([t.r.m.abdelaal-1@tudelft.nl](mailto:t.r.m.abdelaal-1@tudelft.nl))

Lieke Michielsen<sup>1,2#</sup> ([l.c.m.michielsen@student.tudelft.nl](mailto:l.c.m.michielsen@student.tudelft.nl))

Davy Cats<sup>3</sup> ([d.cats@lumc.nl](mailto:d.cats@lumc.nl))

Dylan Hoogduin<sup>3</sup> ([ddhoogduin@gmail.com](mailto:ddhoogduin@gmail.com))

Hailiang Mei<sup>3</sup> ([H.Mei@lumc.nl](mailto:H.Mei@lumc.nl))

Marcel J.T. Reinders<sup>1,2</sup> ([m.j.t.reinders@tudelft.nl](mailto:m.j.t.reinders@tudelft.nl))

Ahmed Mahfouz<sup>1,2\*</sup> ([a.mahfouz@lumc.nl](mailto:a.mahfouz@lumc.nl))

<sup>1</sup> Leiden Computational Biology Center, Leiden University Medical Center, Einthovenweg 20, 2333ZC, Leiden, The Netherlands

<sup>2</sup> Delft Bioinformatics Lab, Delft University of Technology, Van Mourik Broekmanweg 6, 2628XE, Delft, The Netherlands

<sup>3</sup> Sequencing Analysis Support Core, Department of Biomedical Data Sciences, Einthovenweg 20, 2333ZC, Leiden University Medical Center, Leiden, The Netherlands

# Equal contribution

\* Corresponding author ([a.mahfouz@lumc.nl](mailto:a.mahfouz@lumc.nl))

# Supervised classification



- [Cell type classification] is typically solved by **unsupervised clustering** of cells into groups based on the similarity of their gene expression profiles, followed by cell population annotation by assigning labels to each cluster
- The [cell type] annotation step is cumbersome and **time-consuming** as it involves **manual inspection** of cluster-specific marker genes
- **Manual annotations**, which are often not based on standardized ontologies of cell labels, **are not reproducible** across different experiments within and across research groups

# Supervised classification



- In this study, we evaluated the performance of 20 different methods for automatic cell identification using eight scRNA-seq datasets. Several classifiers accurately performed on almost all datasets, particularly:  
**scPred, SVM, scmapcell/cluster, singleCellNet, scVI, LDA and ACTINN.**

# scPred: Cell type prediction at single-cell resolution

José Alquicira-Hernández<sup>1,2</sup>, Anuja Sathe<sup>3,4</sup>, Hanlee P Ji<sup>3,4</sup>, Quan Nguyen<sup>1</sup>, and Joseph E Powell<sup>1,2,5</sup>

<sup>1</sup>Institute for Molecular Bioscience, University of Queensland, Brisbane

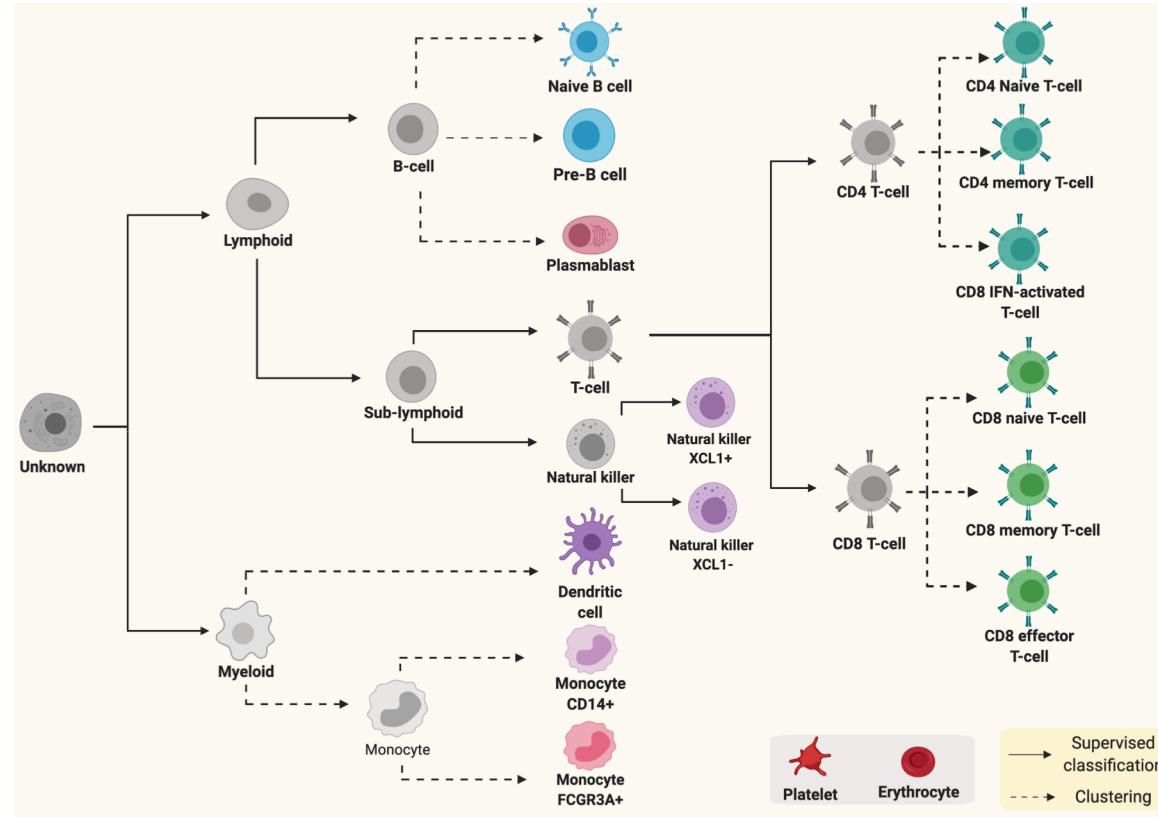
<sup>2</sup>Garvan Weizmann Centre for Cellular Genomics, Garvan Institute of Medical Research, Darlinghurst, Sydney

<sup>3</sup>Division of Oncology, Department of Medicine, Stanford University School of Medicine, Stanford, United States

<sup>4</sup>Stanford Genome Technology Center, Stanford University, Palo Alto, United States

<sup>5</sup>Faculty of Medicine, University of New South Wales, Darlinghurst, Sydney

# Hierarchical classification of PBMCs using supervised and unsupervised approaches



# Recommendations



- **Biology, biology, and biology!**
  - Know your population of study
  - How many cell types do you expect?
  - Expected cell proportions by cell type
  - Expected cell proportions by individual
  - Outlier population or cell contamination? E.g. erythrocytes, platelets
  - If a cell hierarchy is known, use it to guide the classification
  - Identify factors that can confound the prediction
    - Sequencing depth
    - Pool/batch
  - Pay attention to feature selection
- ...

# Recommendations



- Avoid *cluster-hacking*:

“**Exhaustive** clustering search so your data matches your hypothesis”

# Acknowledgements



- Quan Nguyen
- Anuja Sathe
- Hanlee P Ji
- Joseph Powell
  
- Anne Senabouth
- Seyhan Yazar
- Drew Neavin
- Venessa Chin
- Angela Murphy

Single Cell and Computational Genomics

Garvan Institute of Medical Research

Institute of Molecular Bioscience

University of Queensland