# From biology to statistics, and back
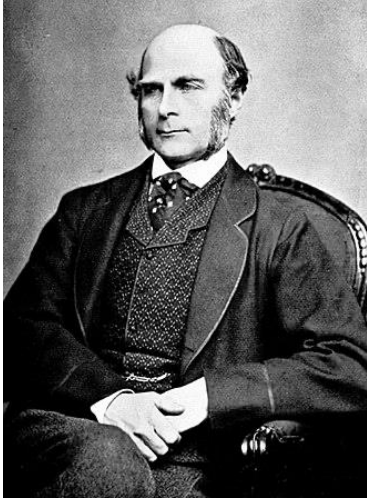
Dr Jiadong Mao

Melbourne Integrative Genomics
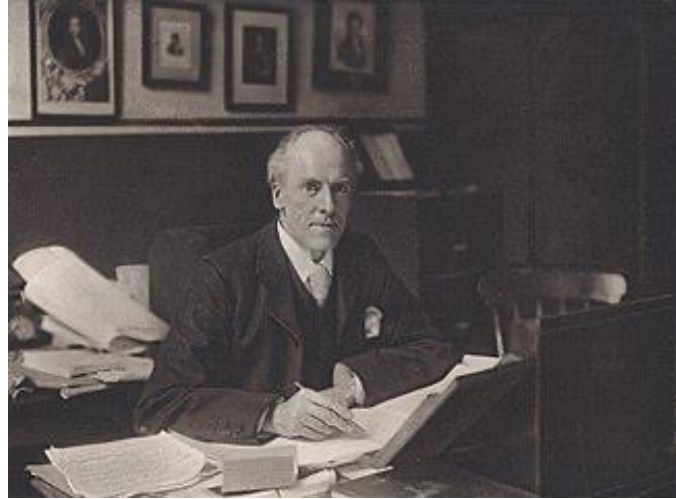
School of Mathematics and Statistics

University of Melbourne

# Statistics and biology



Francis Galton (1822–1911)



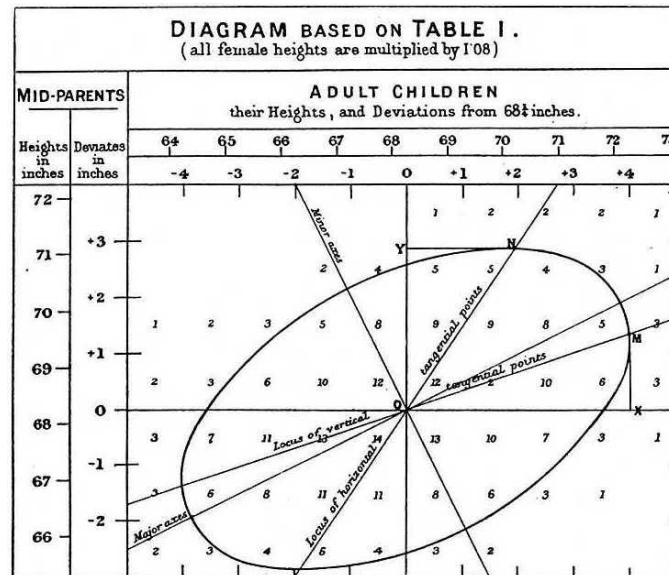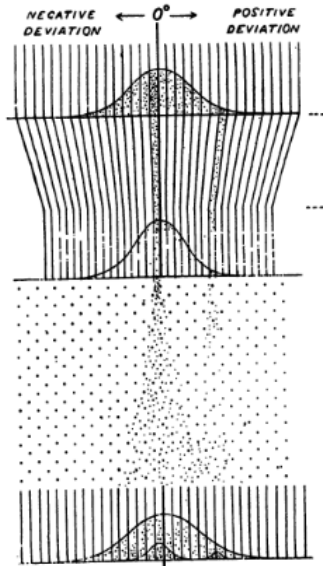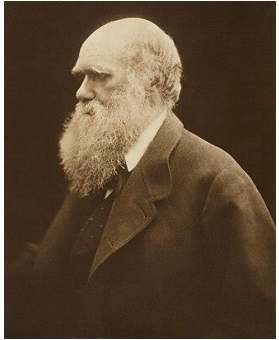Karl Pearson (1857–1936)



Ronald A Fisher (1890–1962)







Fig credit: Wikipedia

# Heredity: the hidden theme of early statistics


Gregor Mendel (1822–1884)
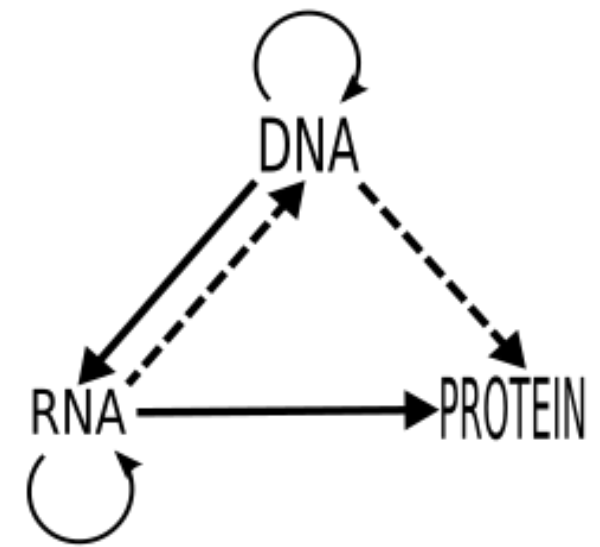

Charles Darwin (1809–1882)

**'Modern synthesis'**
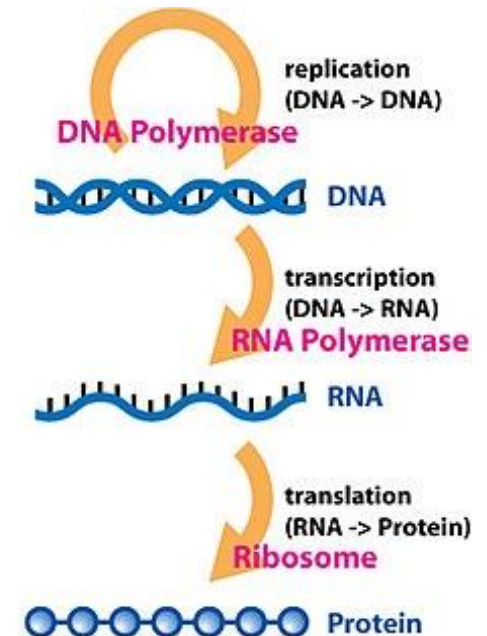

Ronald A Fisher (1890–1962)

(Fisher is) the greatest of Darwin's successors.
-- Richard Dawkins, *The Blind Watch Maker*

Fig credit: Wikipedia

# 'Central dogma' & omics data

- Molecular biology of the cell
- 'Omics' data & high-throughput sequencing
- Types of omics
  - Genomics
  - **Transcriptomics**
  - Proteomics
  - Metabolomics, epigenomics, …
- What's so special about RNAs (transcripts/gene expression)
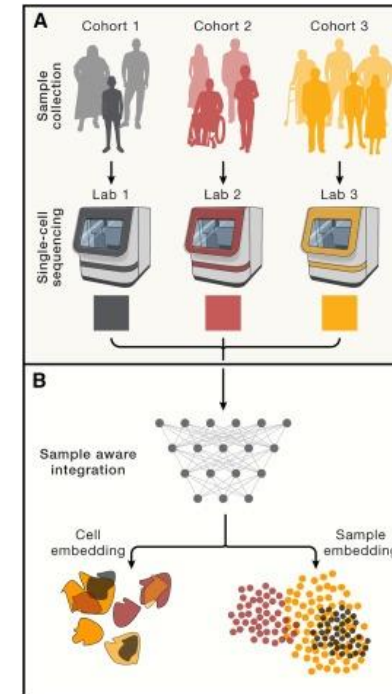  - Cell activities and identities: T cell, B cell, …
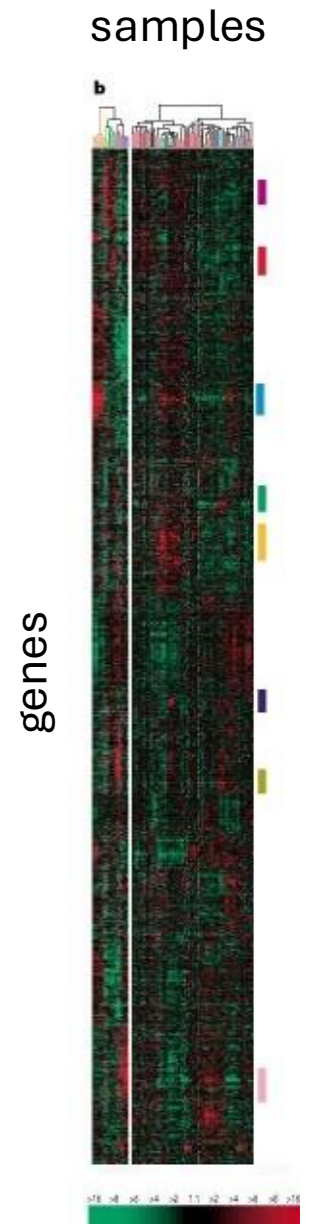


Crick's central dogma (Wikipedia)



Watson's central dogma (Wikipedia)

# Omics and modern statistics

- Bulk RNA sequencing, eg microarray
- HDLSS: high dimension, low sample size
  - ~50 samples, >10,000 genes
- Variable selection, multiple testing
- Common goal: marker gene identification
  - Diagnosis, treatment, prognosis, eg cancer subtyping

samples

genes

Lotfollahi et al. (2024). *Cell*.

Perou et al. (2000). *Nature*.

**Geometric representation of high dimension, low sample size data**

Peter Hall,
*Australian National University, Canberra, Australia*

J. S. Marron
*University of North Carolina, Chapel Hill, USA*

and Amnon Neeman
*Australian National University, Canberra, Australia*

**The Group Lasso for Logistic Regression** FREE
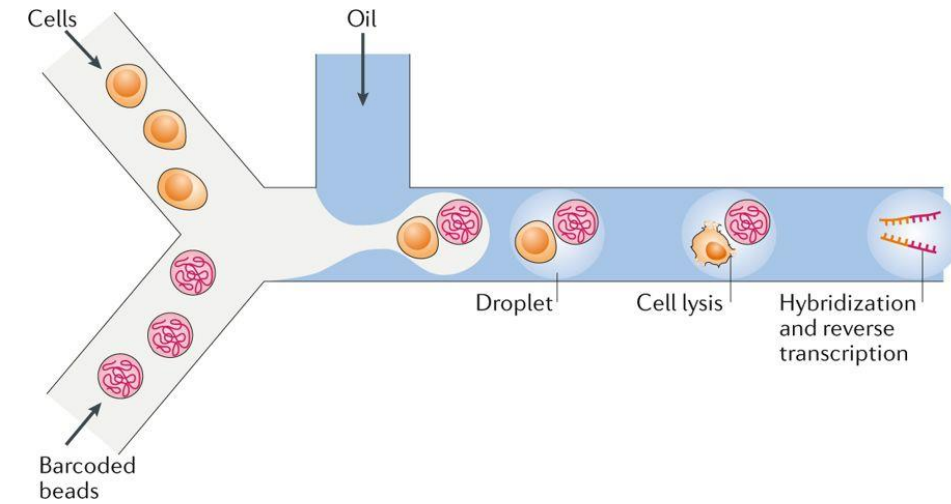
Lukas Meier ✉, Sara Van De Geer, Peter Bühlmann

ON TESTING THE SIGNIFICANCE OF SETS OF GENES

BY BRADLEY EFRON[1] AND ROBERT TIBSHIRANI[2]

*Stanford University*

# New challenges: big omics data



https://www.rna-seqblog.com/wp-content/uploads/2018/08/droplet.png

- Single-cell RNA sequencing (scRNA-seq)
  - Dissolve tissue into single cells & seq

- 'HDHSS': High dimension, high sample size
  - >10,000 cells per sample/donor
  - >20,000 genes

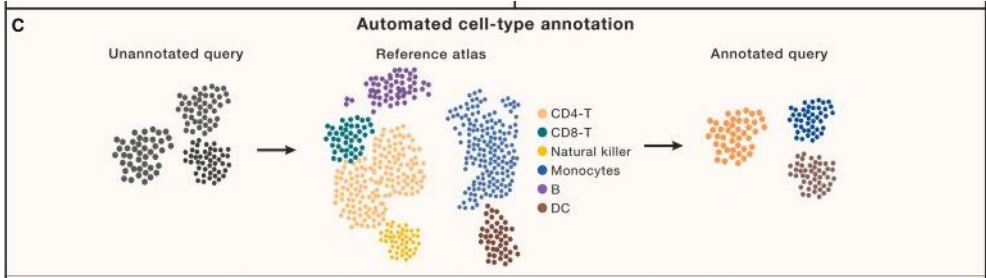- Finding marker genes at cell (type) level
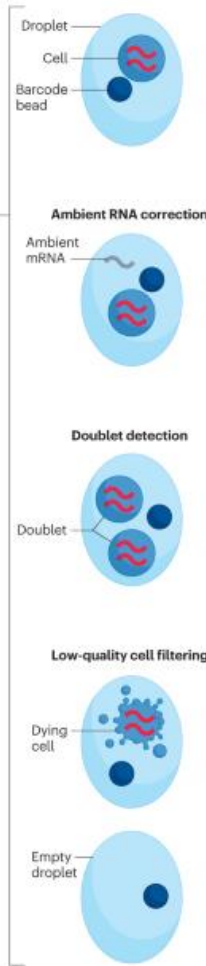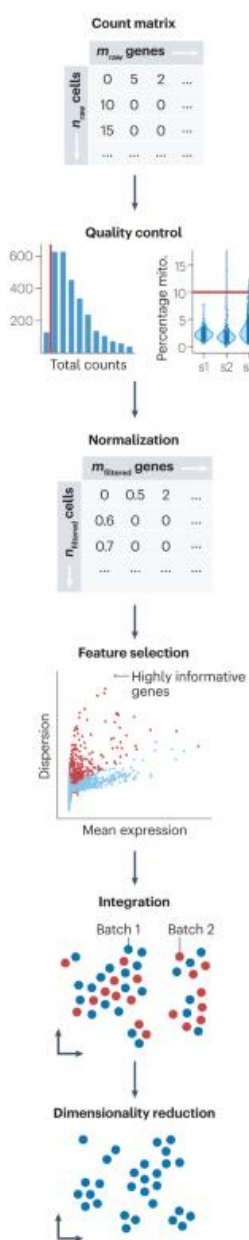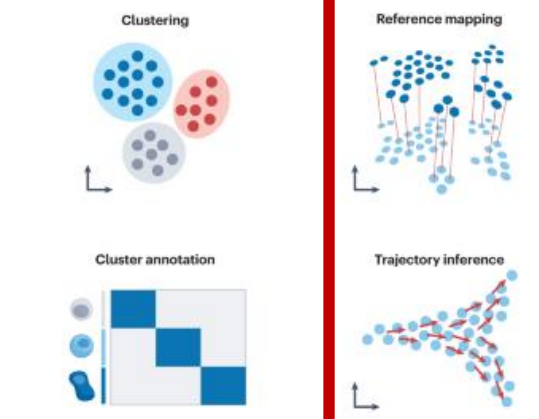


Lo et al. (2020). *Cancers*.

# Cell type annotation

- Cells form (relatively) homogeneous groups
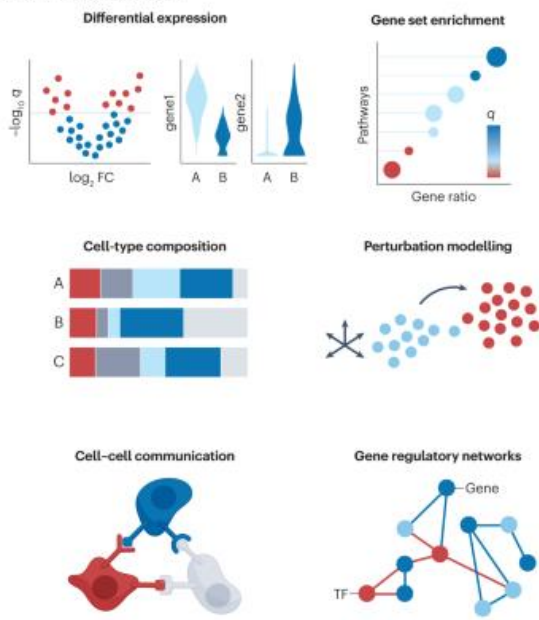- Group cells into cell types: train cell type classifier on reference data



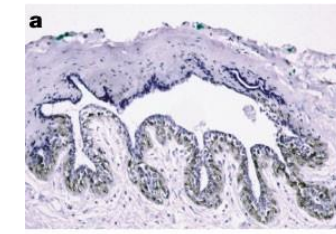Lotfollahi et al. (2024). *Cell*.

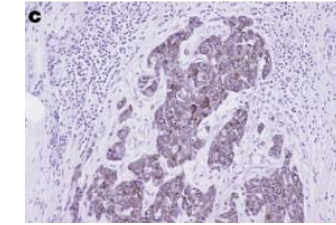Heumos et al. (2023). Nat Review Genetics.

# Adding spatial information

- Spatial transcriptomics (ST)
- Why ST
  - High-throughput: measuring a lot of molecules
  - scRNA-seq: which cells are doing what
  - ST: which cells are doing what, and **where they are doing it**
  - Example: tumour microenvironment
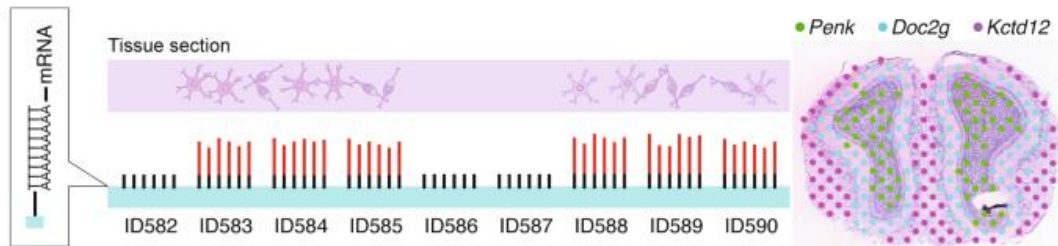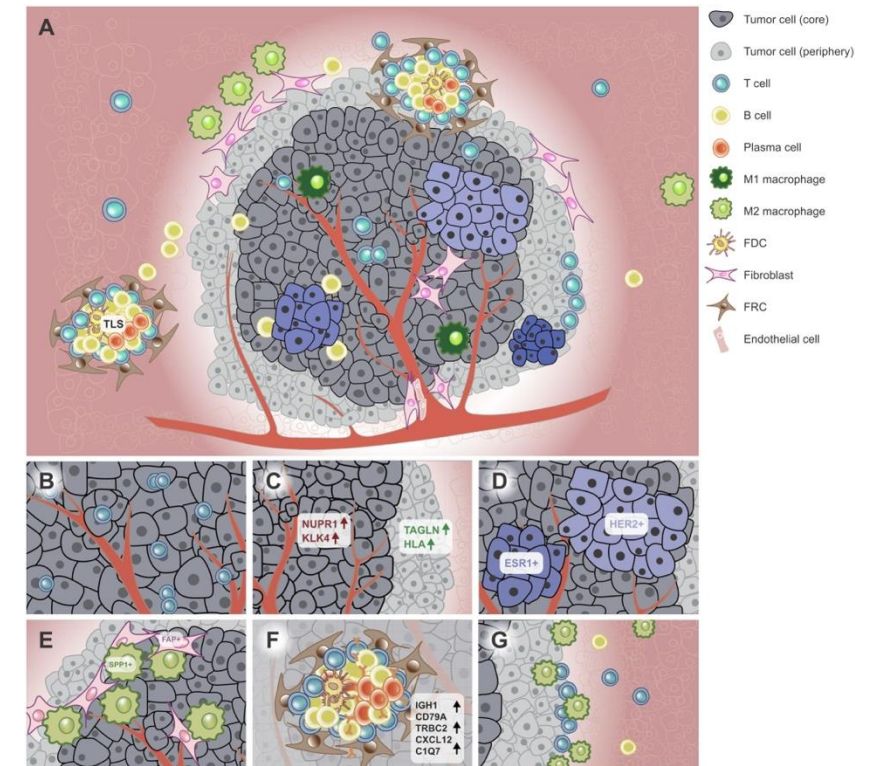- Cell type annotation is still the key



normal

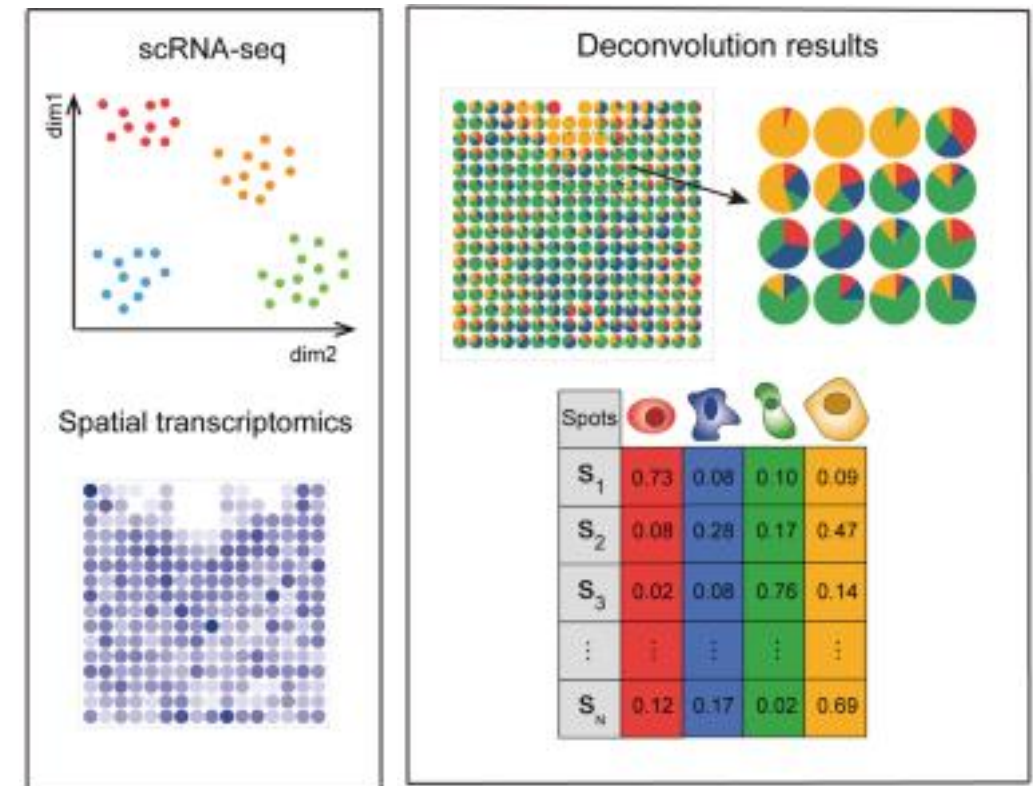tumour

Perou et al. (2000). *Nature*.





Nature Methods Method of the year 2021
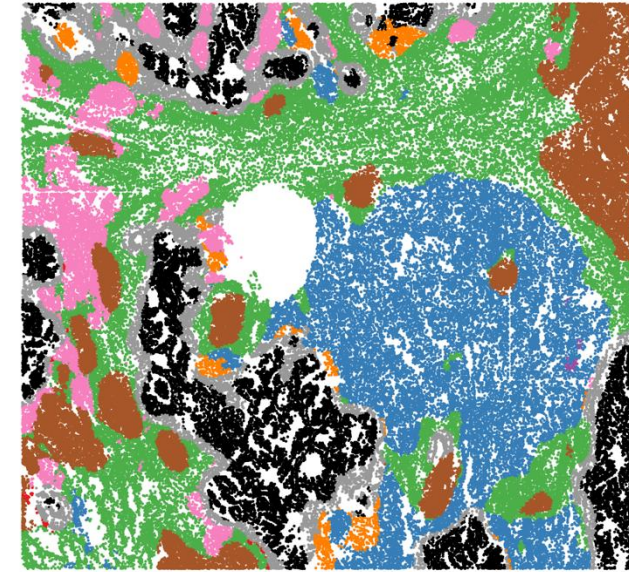
# Cell type deconvolution

- Main idea
    - Each 'spot' may contain multiple cells
    - Use scRNA-seq as 'reference'
    - Decompose gene expression in each spot as combination of reference cell types
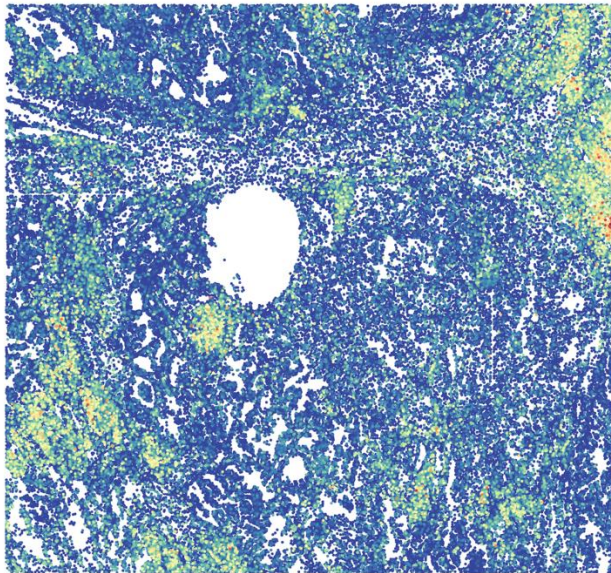


Li et al. (2023). Nat Comms

# Deconvolute cancer cell states

- *Platform*: Nanostring CosMx
- *Sample:* human lung cancer (non-small cell lung cancer)
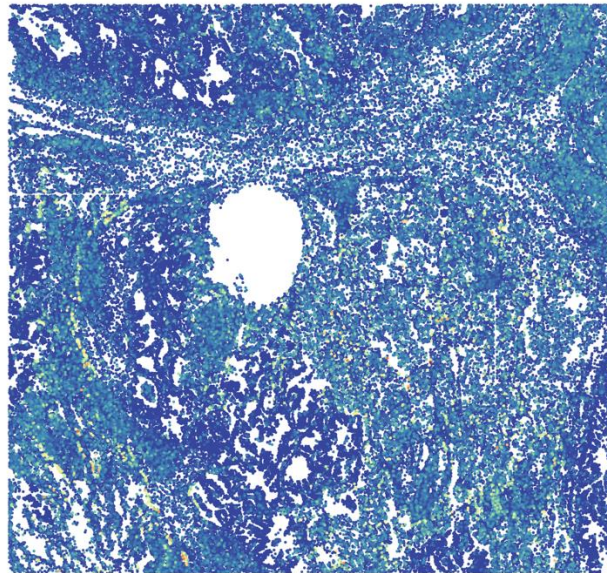- *References*: scRNA-seq from lung fibrosis patients

# ST with cancer cell lineage tracing

Henrietta Holze

Dane Vassiliadis

Mark Dawson

- *Platform*: BGI Stereo-seq
- *Sample info*
  - mouse spleen sample with AML cells
  - SPLINTR lineage tracing, same AML **clone** same **barcode**
- *References*: scRNA-seq from mouse spleen and bone marrow

Peter Mac
Peter MacCallum Cancer Centre
Victoria Australia

Inject barcoded AML cells → AML spleen → Stereo-seq → Gene expression & AML barcode distributions → Annotation & analyses

PhiSpace ST

AML: acute myeloid leukemia
SPLINTR: single-cell profiling and lineage tracing

# Cell states of 'meta-clones'



Density of all AML barcodes

Barcode_25    Barcode_2613

Barcode_6595    Barcode_19949

AML Meta-clones

2  5  7  9
3  6  8  10

**Neutrophil-like**
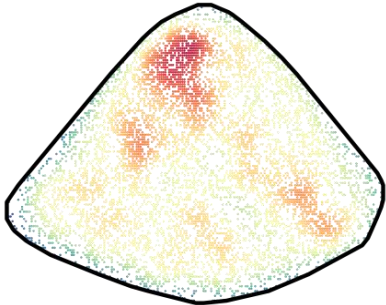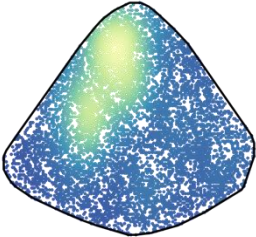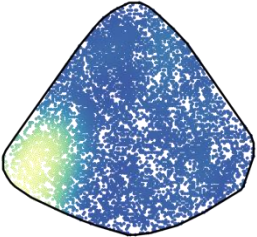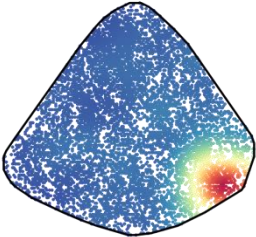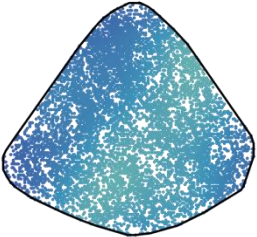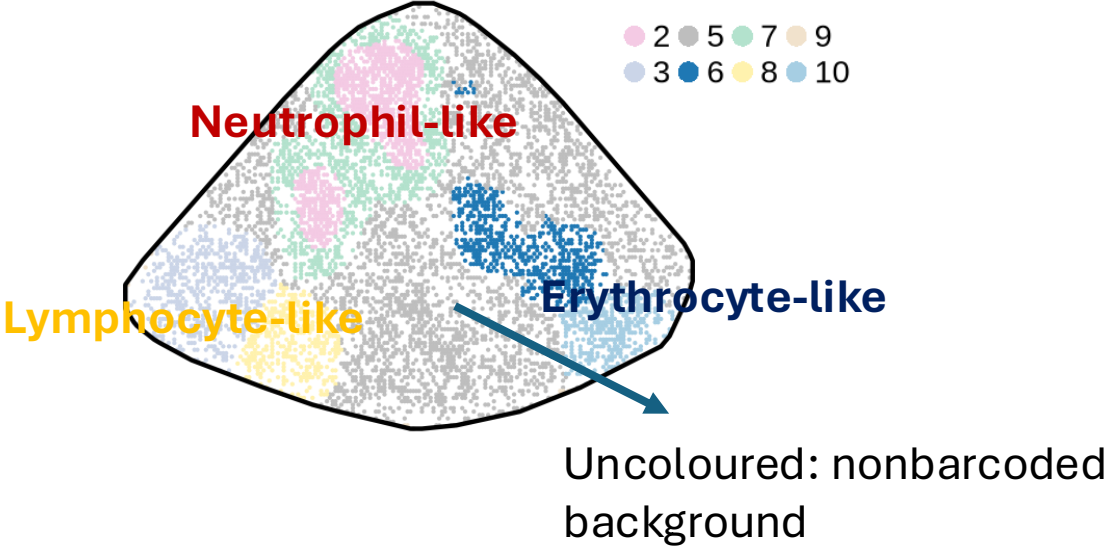
**Lymphocyte-like**    **Erythrocyte-like**

Uncoloured: nonbarcoded background

### Meta-clone-specific enriched cell types

|   | Meta-clone 2 | Meta-clone 3 | Meta-clone 6 | Meta-clone 7 | Meta-clone 8 | Meta-clone 10 |
|---|---|---|---|---|---|---|
| 1 | Granulo (BM) | Imm NK (BM) | RBC (CITE) | HPC (BM) | HPC (BM) | Naive B (BM) |
| 2 | Neutro (Spleen) | Trans B (Spleen) | RBC (Spleen) | IMM 1 (Neutro) | T (BM) | ProEryThBla (BM) |
| 3 | HPC (BM) | T (BM) | ErythBla (BM) | Neutro (Spleen) | Pre-B cycl (Spleen) | Pre-B cycl (Spleen) |
| 4 | T1 (Neutro) | MAT 3 (Neutro) | ProEryThBla (BM) | Pre-B cycl (Spleen) | Imm NK (BM) | ICOS+ Tregs (CITE) |
| 5 | Imm NK (BM) | CD8 T (Spleen) | IMM 2 (Neutro) | Pre-B (BM) | Pre-B (BM) | CD4 T (CITE) |

# Reflections

- Fast-evolving biotech, reliable stats method needed
  - 3D spatial, spatio-temporal, ...
- Collaborative culture
  - Wet: Biologists, bioinformaticians;
  - Bridge: computational biologists;
  - Dry: statisticians, mathematicians, computer scientists
- (Effective) visualisation
  - Most commonly used: ggplot2 & plotly
  - What you want to show ≠ what viewers see