# AI LAB

## 2.1 - UNSUPERVISED LEARNING

FRANCESCA M. BUFFA

# LAB STRUCTURE

- INTRODUCTION
- THE DATA
- THE AI-LAB CHALLENGE – PART 1
- PART 1 - SHARING AND DISCUSSION

- UNSUPERVISED LEARNING EXAMPLES
- THE AI-LAB CHALLENGE – PART 2
- PART 2 - SHARING AND DISCUSSION

- SUPERVISED LEARNING EXAMPLES
- THE AI-LAB CHALLENGE – PART 3
- LARGE PROJECTS AND DATABASES

- THE AI-LAB CHALLENGE PARTS 1-3, SHARING AND DISCUSSION
- DATA INTERPRETATION
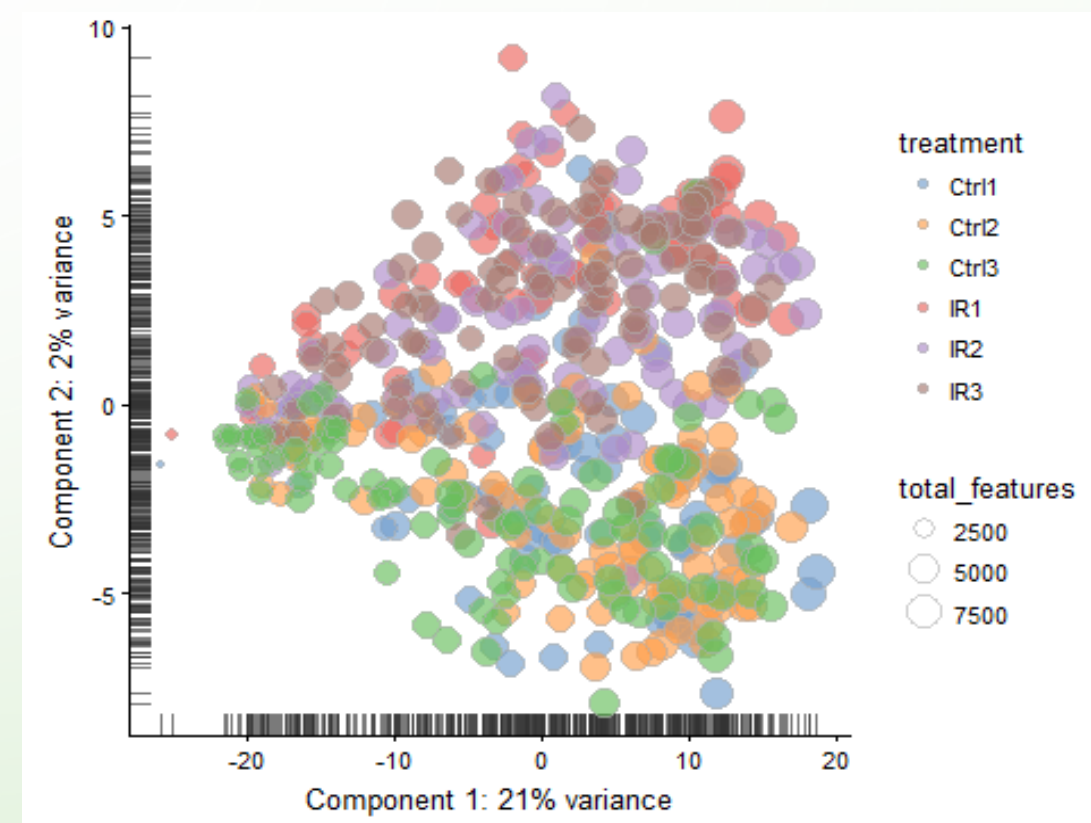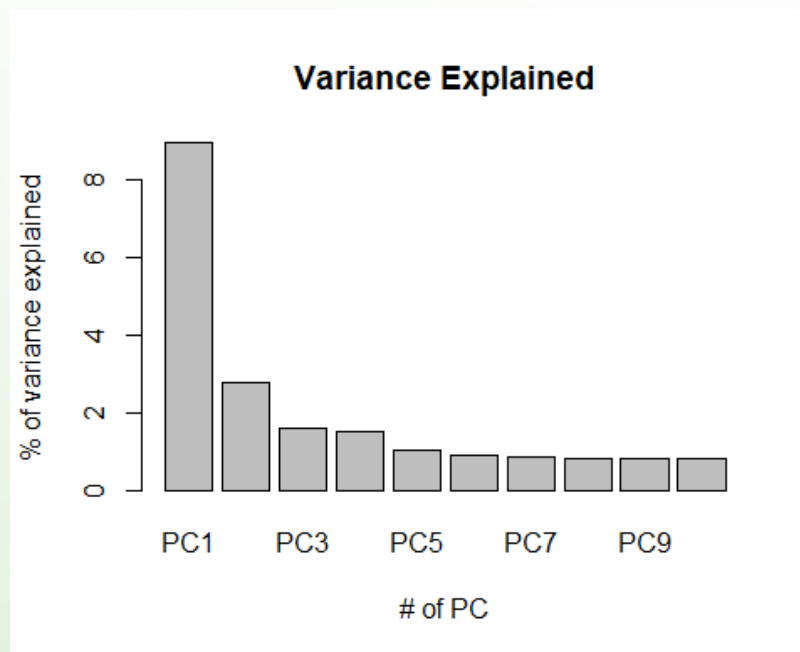- DISCUSS AND PREPARE WORKSHOP PRESENTATIONS

# UNSUPERVISED LEARNING APPLICATIONS

- EXPLORATORY DATA ANALYSIS

- QUALITY CONTROL

- CLASS DISCOVERY

- DIMENSIONALITY REDUCTION
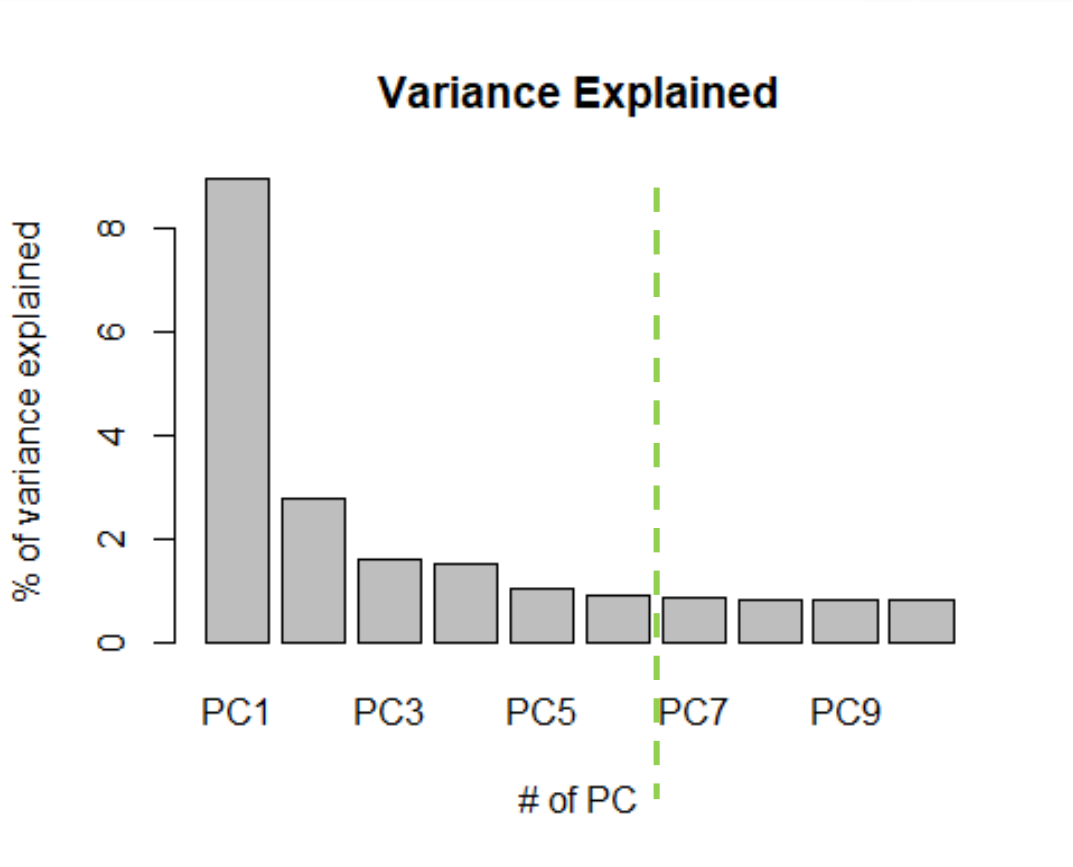
# EXAMPLE: QC IN SINGLE CELL SEQUENCING EXPERIMENT
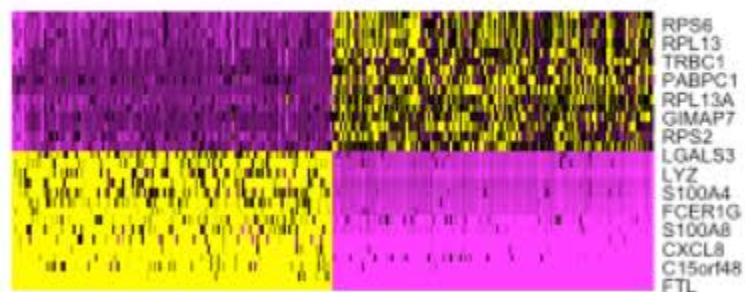
# PCA DIMENSIONALITY REDUCTION

**#Cells 528**
**# Samples (Cntrl vs Treatment)**
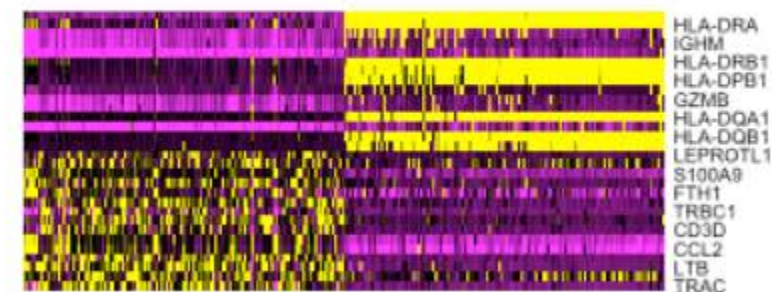
**#Features/Genes 2637**

# PCA DIMENSIONALITY REDUCTION

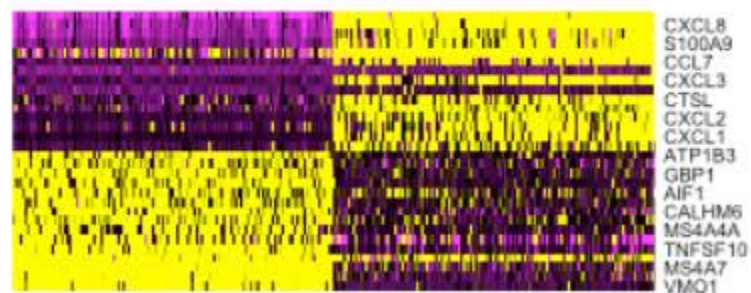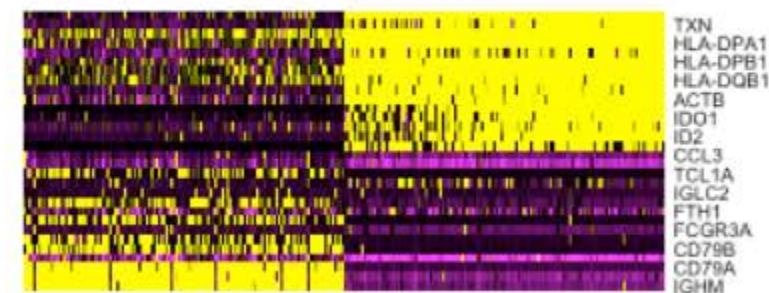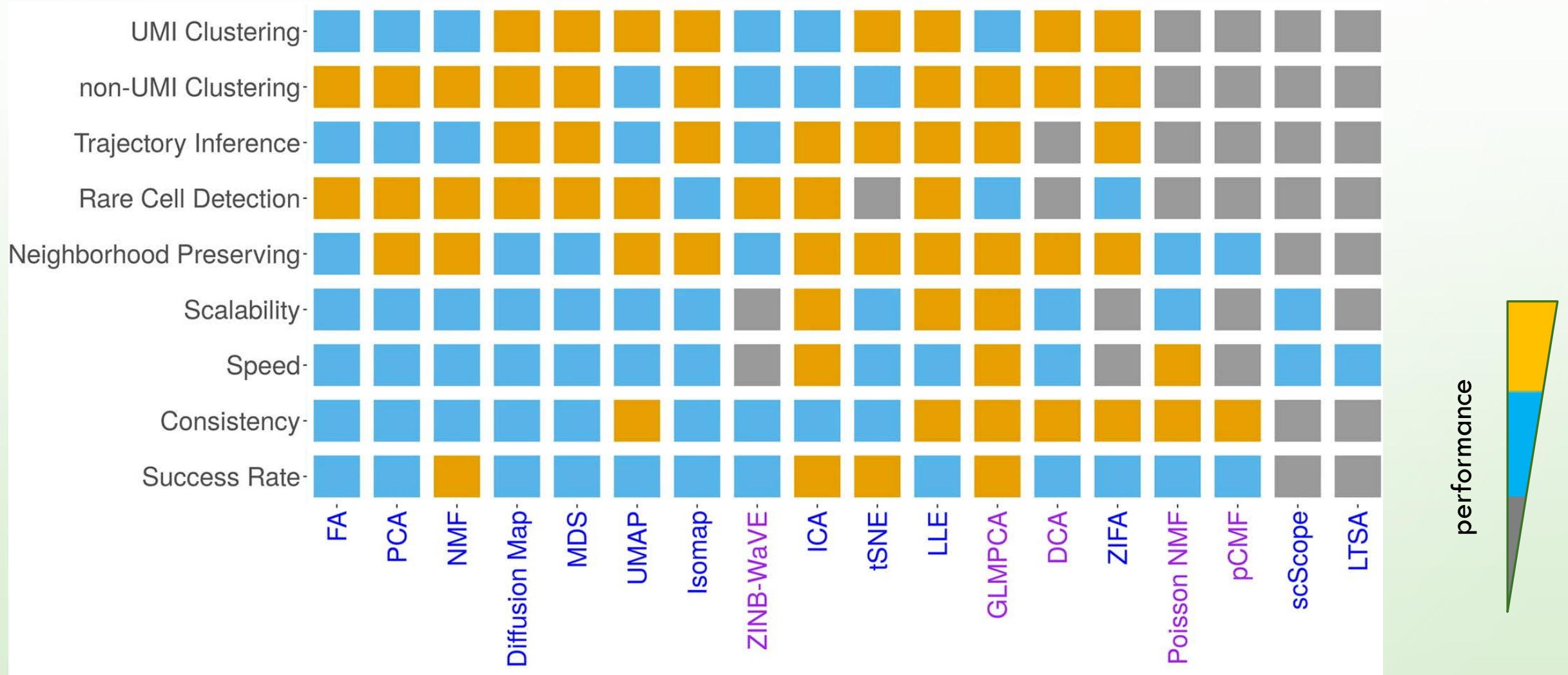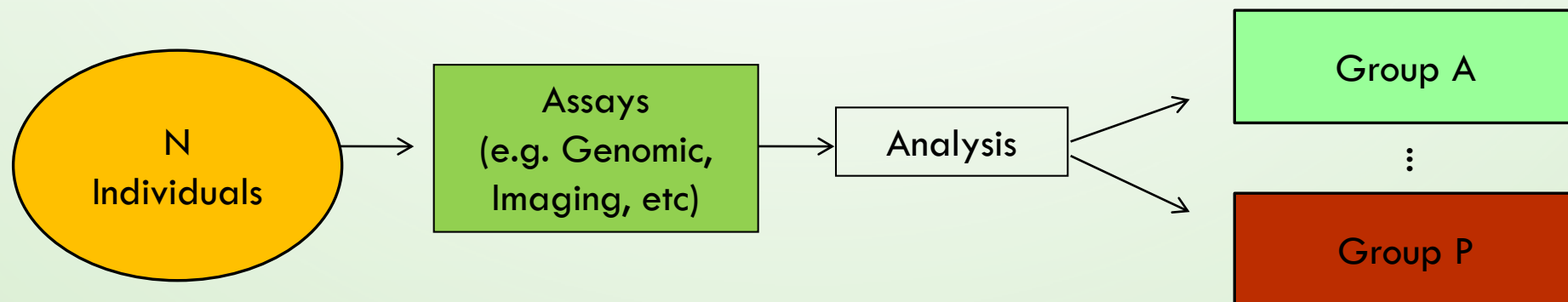# DIMENSIONALITY REDUCTION

# CLUSTERING APPLICATIONS

## ARE THERE GROUPS OF SIMILAR DISEASES/CELLS?

"UNSUPERVISED APPROACH"

CLASS DISCOVERY

FIND GROUPS OF SIMILAR CASES OR SIMILAR FEATURES
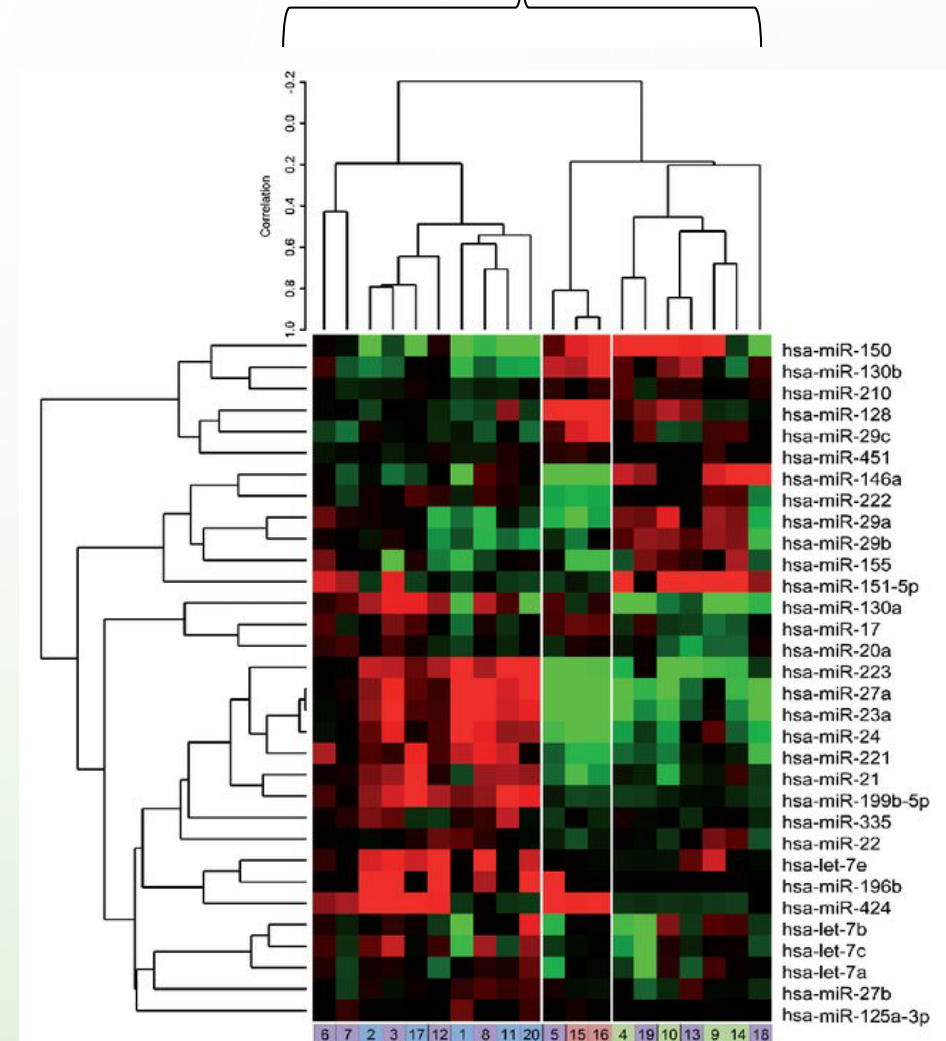
# Clustering: guilty-by-association

**Are my samples similar?**
Samples with similar genomic profile might for example have a similar prognosis or response to treatment.
Or in single cell might come from the same cell population.
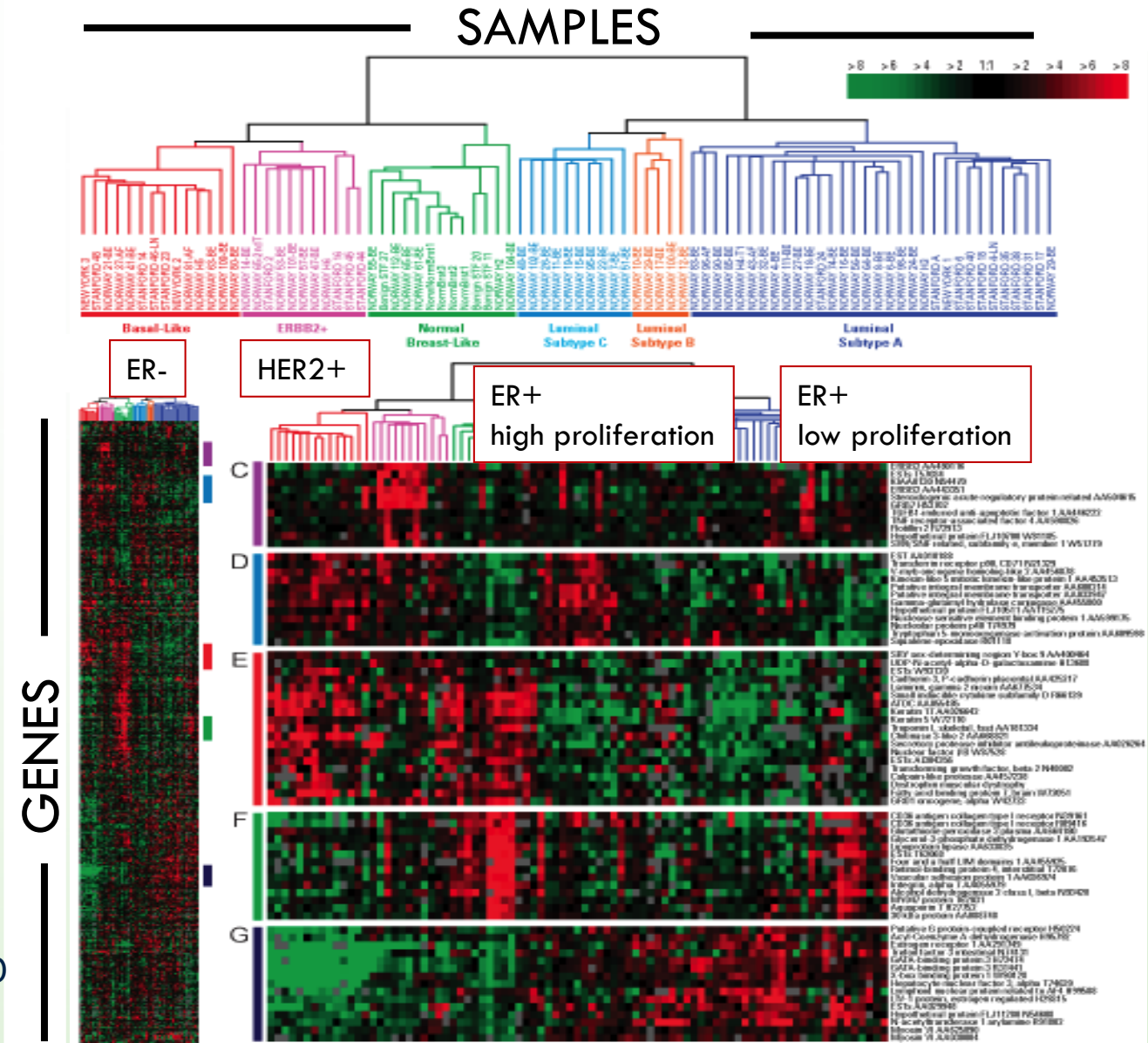
**Are my genes similar?**
Suppose genes A and B are grouped in the same cluster. This mean they are expressed under the same conditions.
Then we can hypothesize that genes A and B are involved in similar pathways/share function.
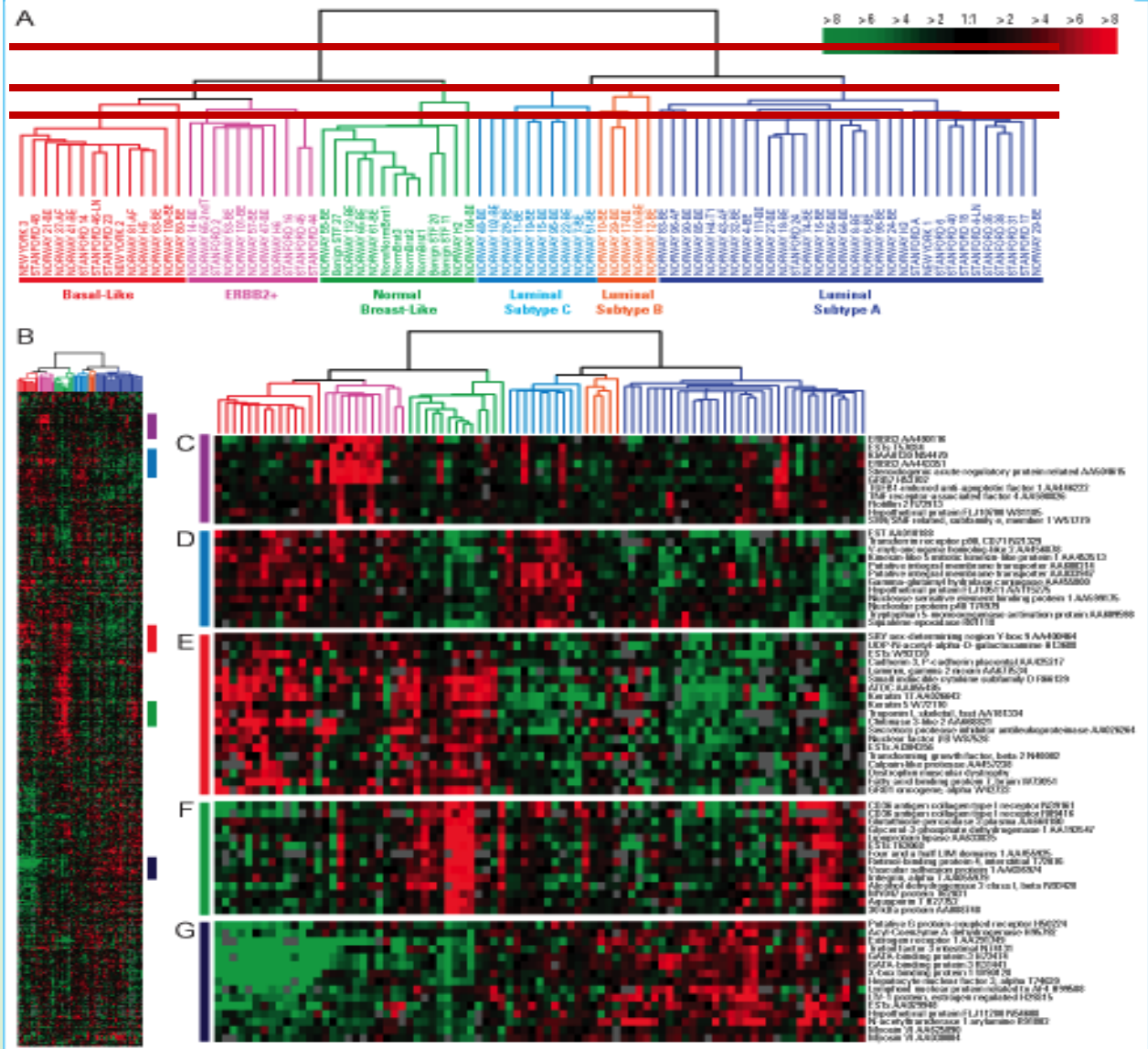
# Unsupervised: Breast Cancer Subtypes (PAM50)

- Unsupervised approach: outcome not considered

- Classification based on Expression of Breast Cancer Intrinsic Genes

- Gene expression microarrays

- Hierarchical clustering used to represent distance between samples

- Groups identified matching existing clinical knowledge

Perou et al, Nature 2000
Sorlie et al, PNAS 2001

# How many clusters?



Perou et al, Nature 2000
Sorlie et al, PNAS 2001
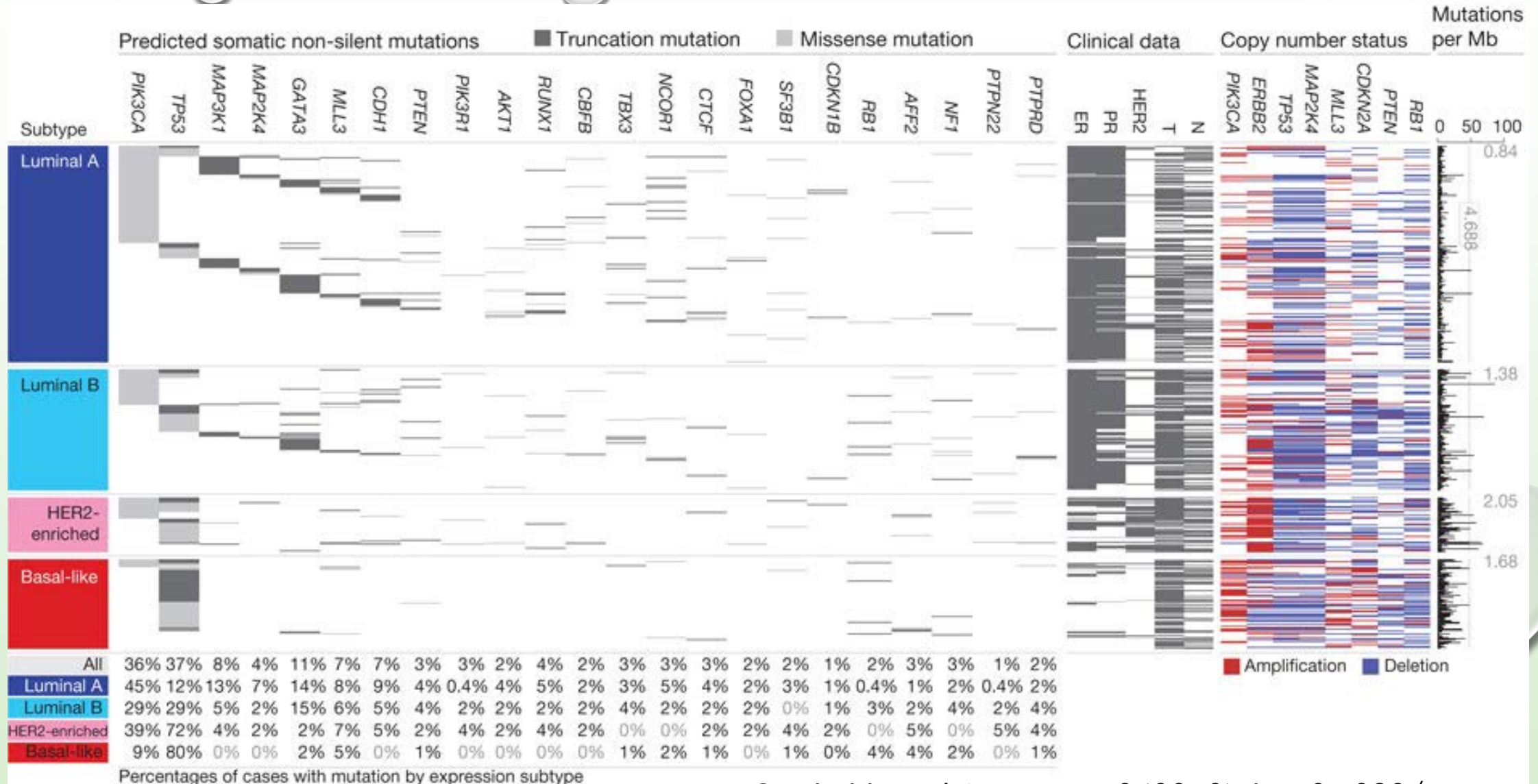
# REPRODUCIBILITY

WEIGELT, B ET AL, LANCET ONCOLOGY, 2010

WHEN USING DIFFERENT METHODS TO ASSIGN PATIENTS TO EACH CLUSTER:
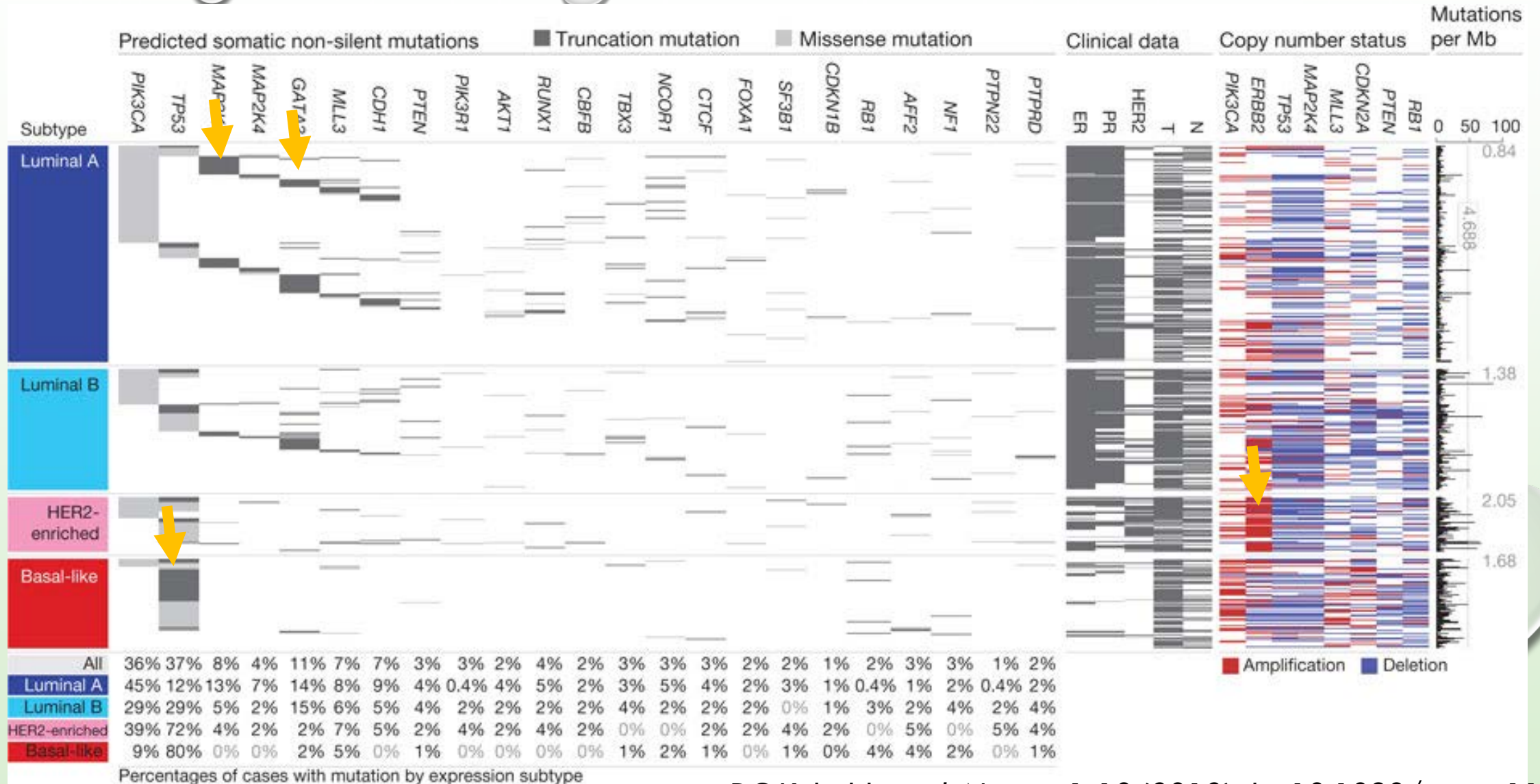
- BASAL-LIKE CANCERS WERE CONSISTENTLY CLASSIFIED.

- ASSIGNMENT OF INDIVIDUAL CASES TO LUMINAL A, LUMINAL B, HER2, AND NORMAL BREAST-LIKE SUBTYPES WAS DEPENDENT ON THE METHOD USED.

- THE SIGNIFICANCE OF ASSOCIATIONS WITH OUTCOME OF EACH MOLECULAR SUBTYPE, OTHER THAN BASAL-LIKE AND LUMINAL A, VARIED DEPENDING ON THE METHOD USED.

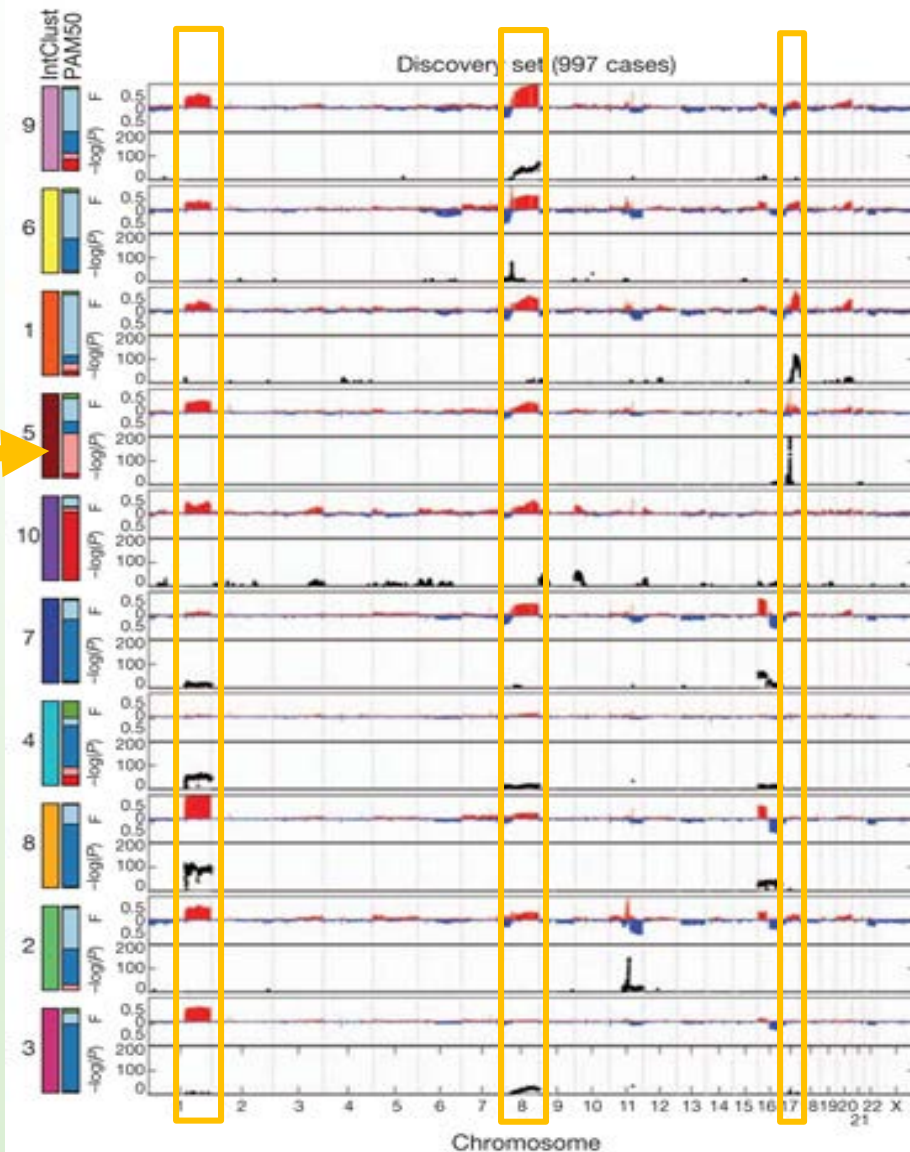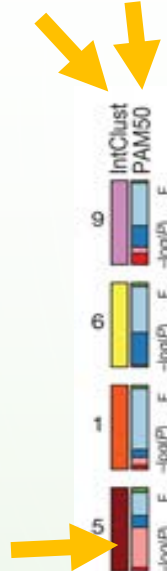# Link between gene expression clusters and genomic features

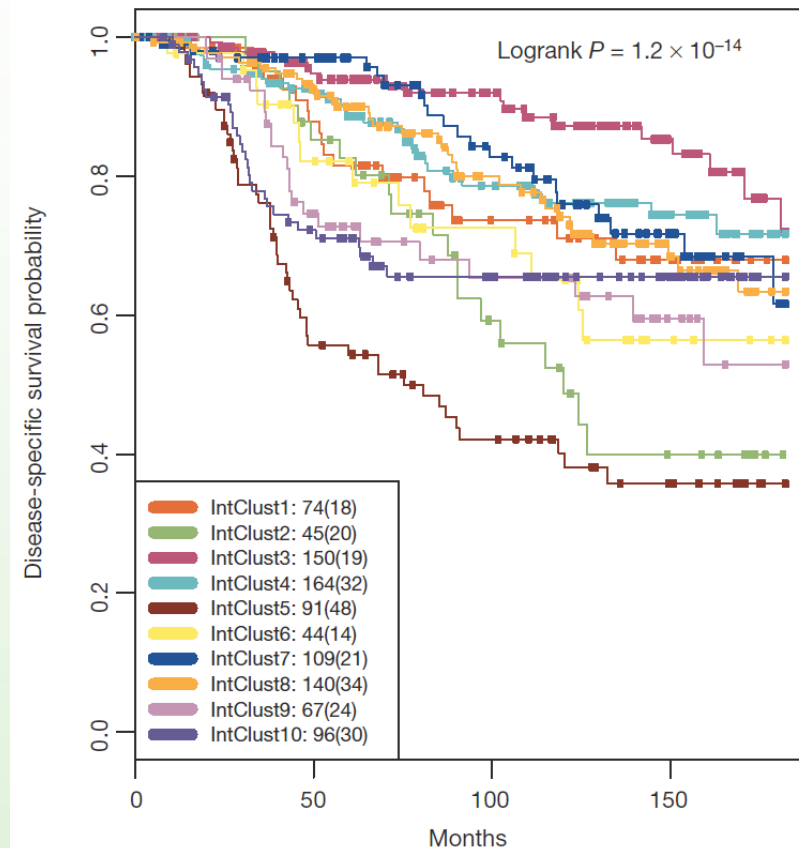# Link between gene expression clusters and genomic features

# Clustering on Copy Number Changes and transcriptional landscape of thousands of tumours (*IntClust*)



The genomic and transcriptomic architecture of 2,000 breast tumours reveals novel subgroups

# Gene expression-based approach for classifying breast tumors into the ten IntClust subtypes
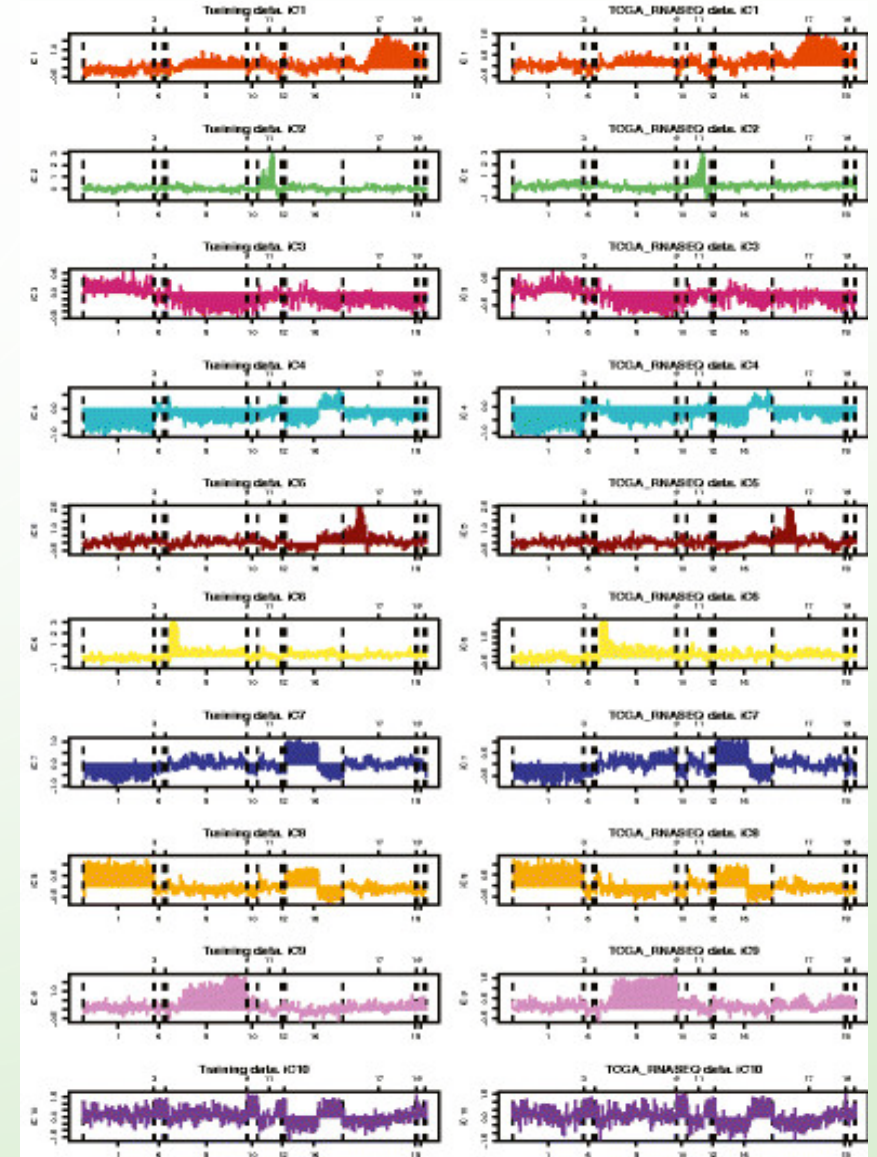


Genome-driven integrated classification of breast cancer validated in over 7,500 samples

H Raza Ali, Oscar M Rueda, Suet-Feung Chin, Christina Curtis, Mark J Dunning, Samuel AJR Aparicio and Carlos Caldas ✉

*Genome Biology* 2014 **15**:431 | DOI: 10.1186/s13059-014-0431-1 | © Ali et al.; licensee BioMed Central Ltd. 2014
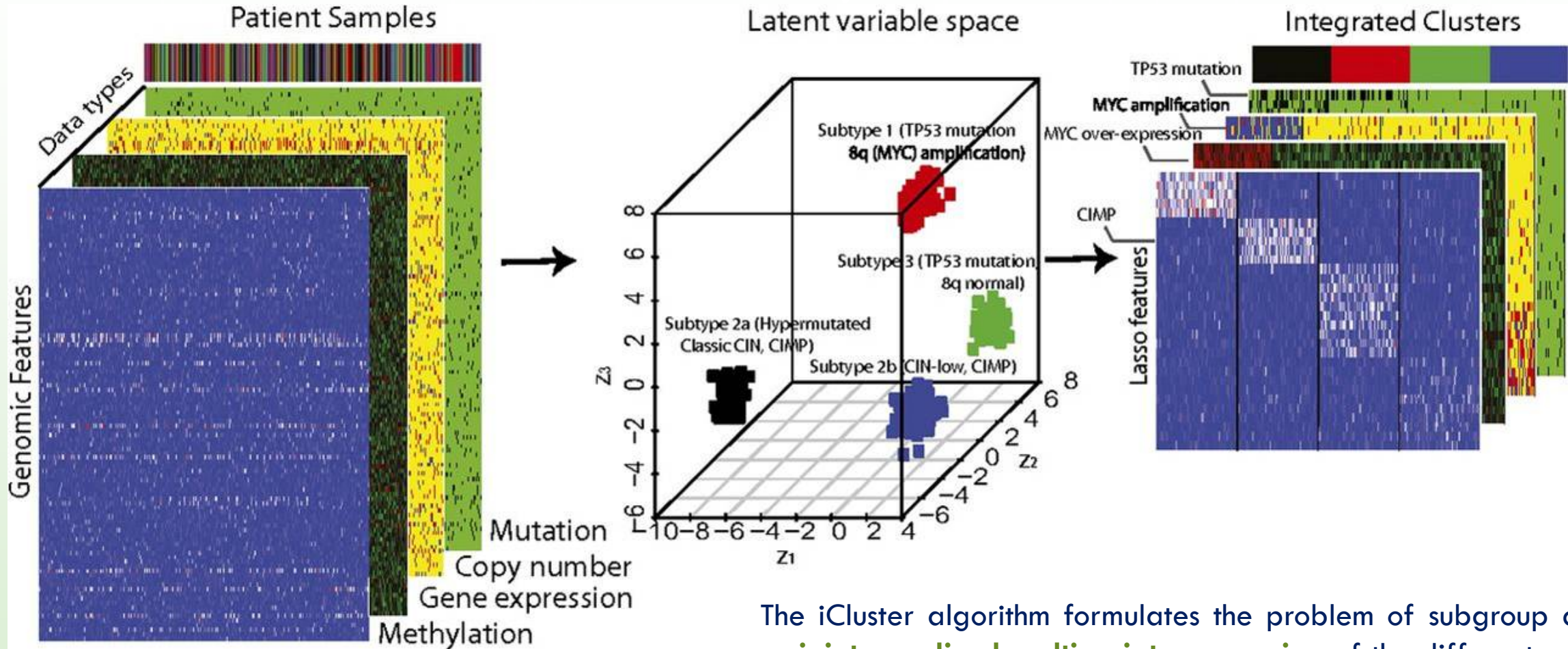
# EXAMPLE OF INTEGRATED CLUSTERING METHOD: ICLUSTER
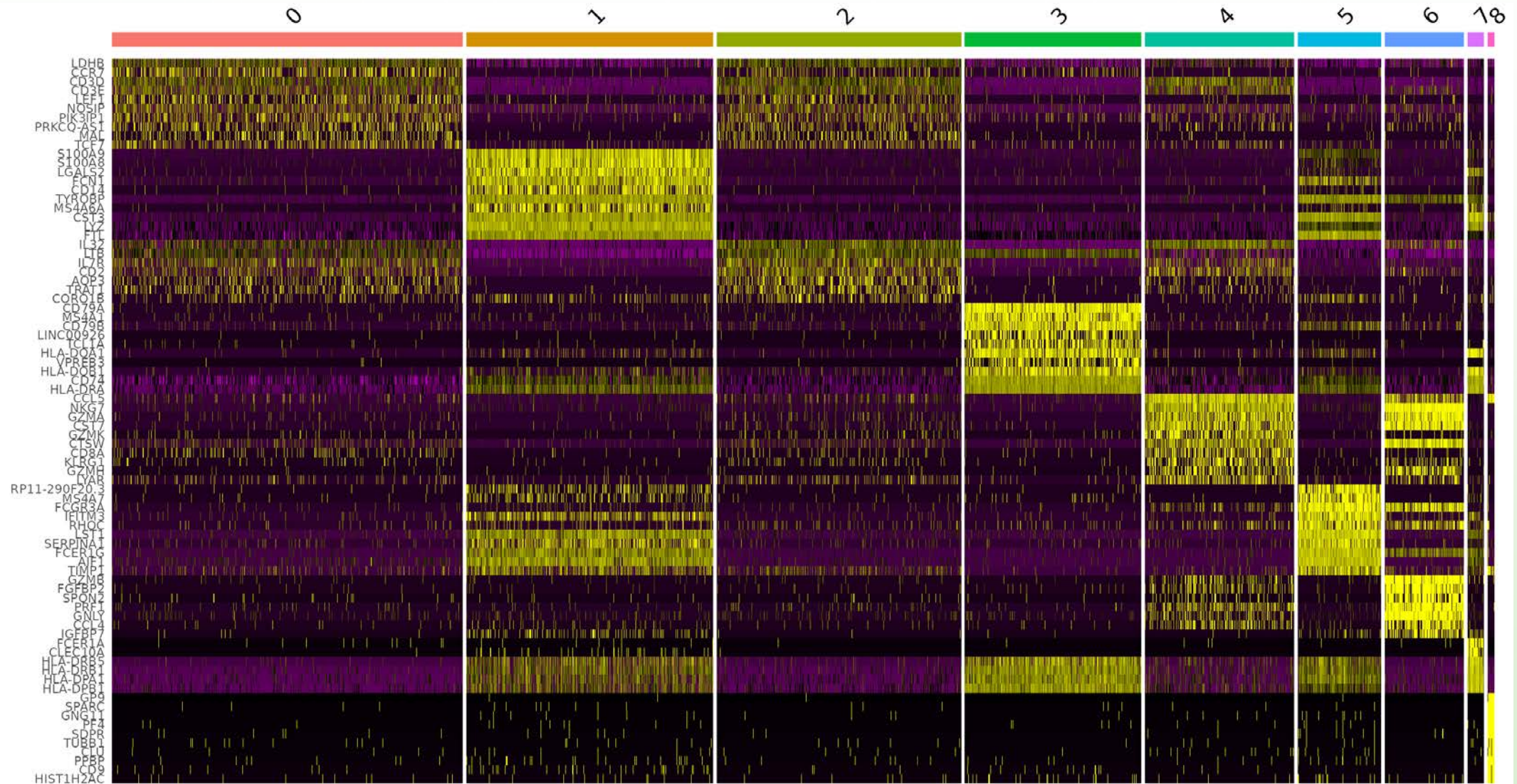


*https://sites.google.com/site/ronglais/icluster*

The iCluster algorithm formulates the problem of subgroup discovery as **joint penalized multivariate regression** of the different omics data types with **reference to a set of common latent variables**, which represent the underlying tumor subtypes Gaussian joint latent variable model.
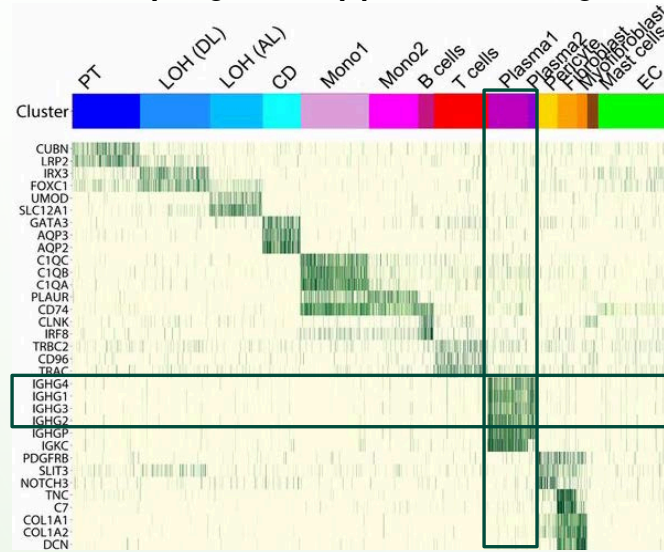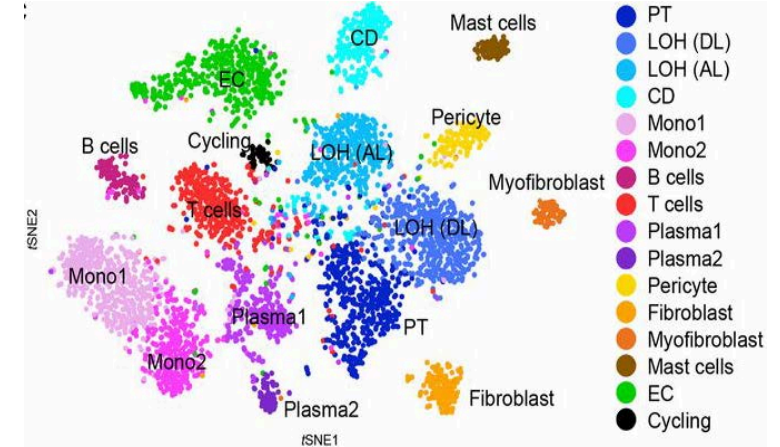
# CLUSTERING SINGLE CELLS RNASEQ

# DISCOVERY OF DIFFERENT CELL PHENOTYPES

Identifying cell-type marker genes

Discovering sample heterogeneity



Uncovering tissue dynamics