# Data integration in cancer research

## An overview of the existing approaches
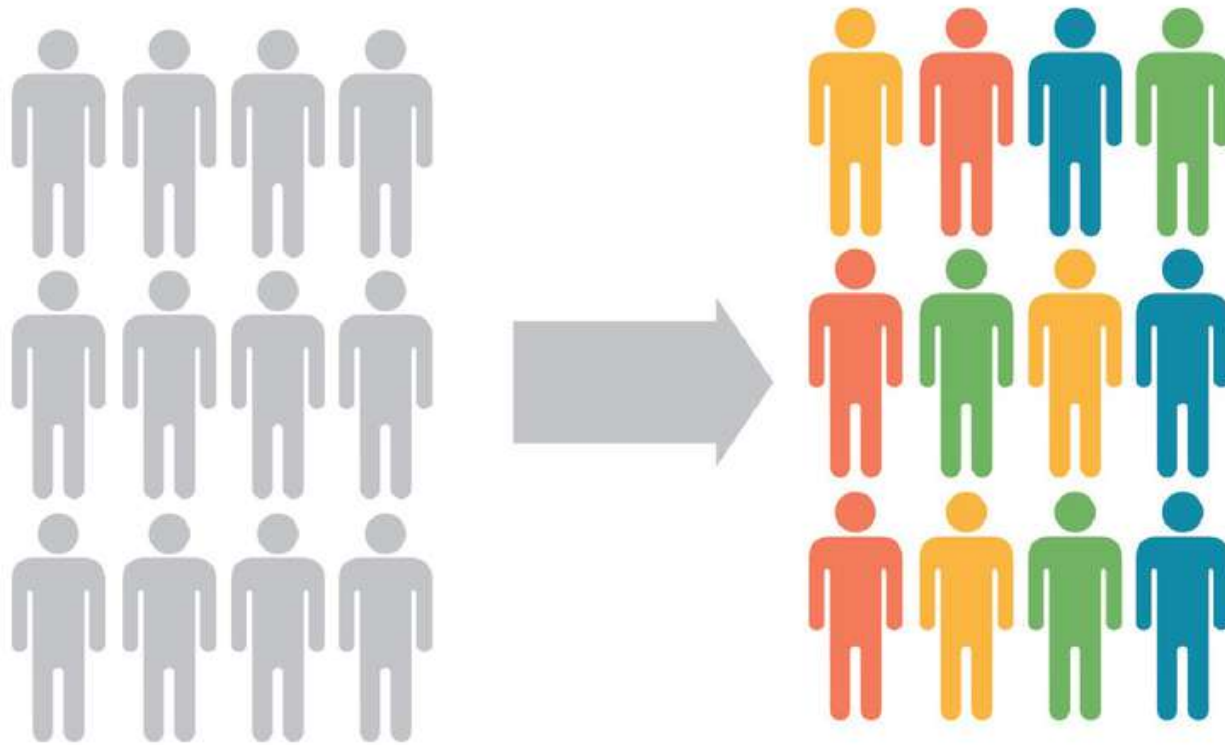
**Laura Cantini**
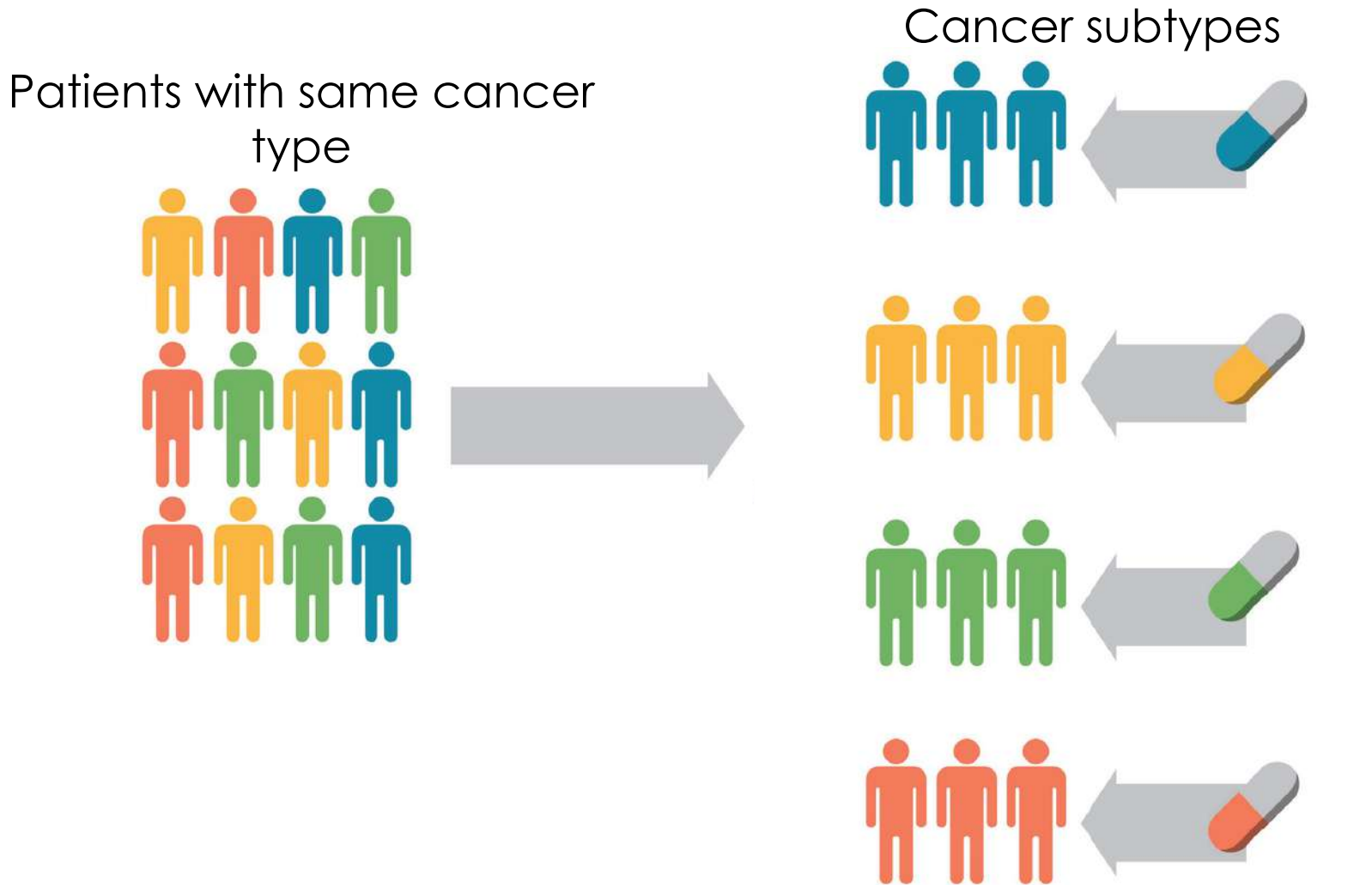Computational Systems
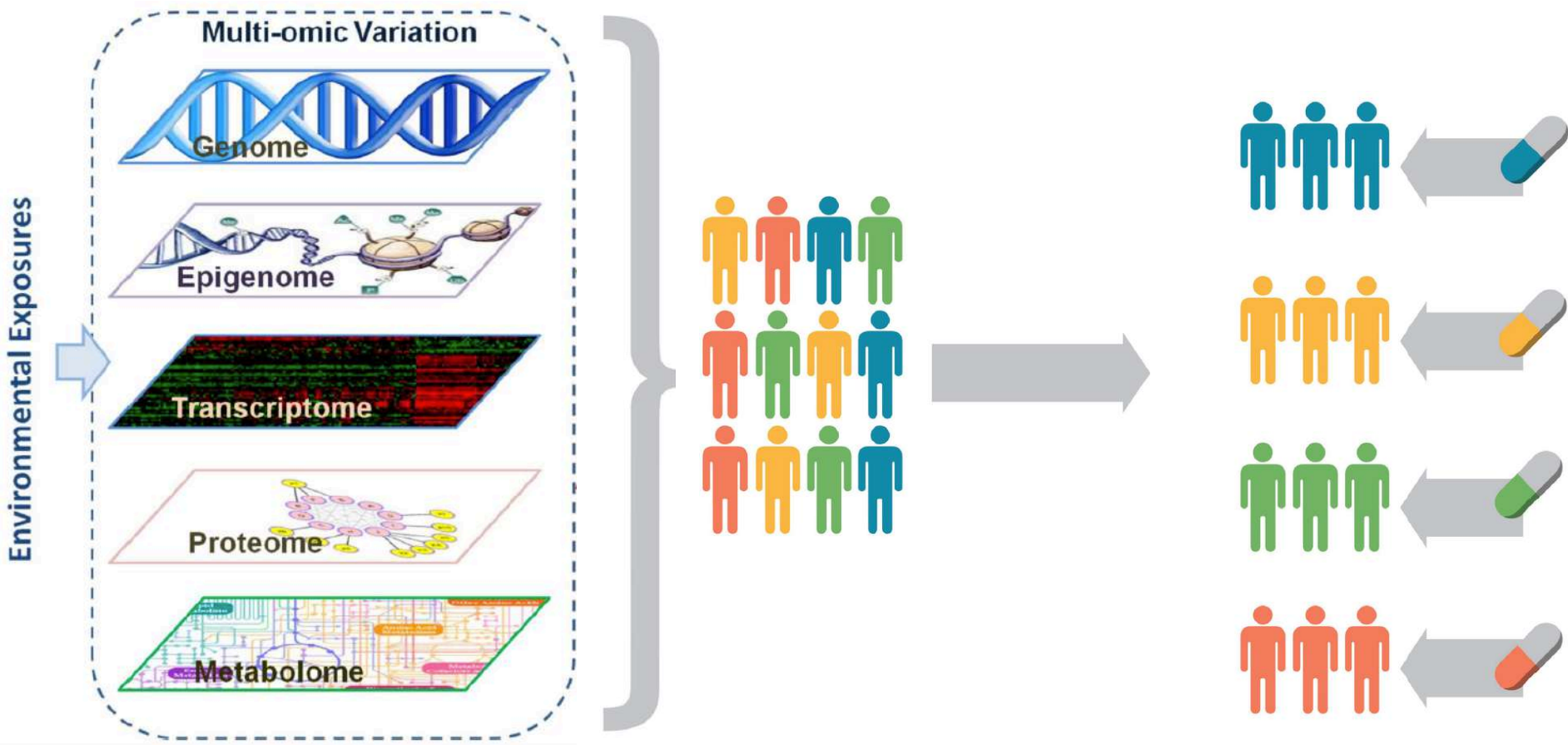Biology Team
IBENS, Paris

# Personalized cancer medicine

Patients with same cancer type don't have the same survival, treatment response and molecular characteristics
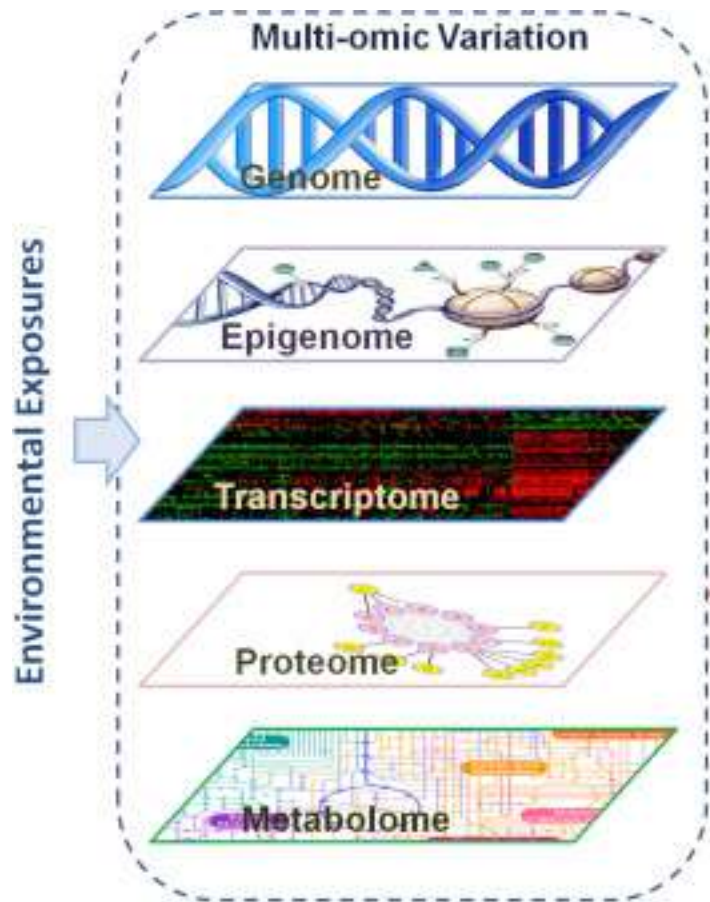
# Personalized cancer medicine



Patients with same cancer type

Cancer subtypes

Classify cancer patients into groups with similar prognosis, drug response or molecular features

# Multi-omics data available



**The Cancer Genome Atlas (TCGA)** for example contains data from 10.000 patients, 33 cancer types, 6 omics, plus clinical data

The Cancer Genome Atlas Research Network, Weinstein, J.N., Collisson, E.A., Mills, G.B., Shaw, K.M., Ozenberger, B.A., Ellrott, K., Shmulevich, I., Sander, C., and Stuart, J.M. (2013) The Cancer Genome Atlas Pan-Cancer analysis project. Nat Genet. doi:10.1038/ng.2764

Sun, Yan V., and Yi-Juan Hu. "Integrative analysis of multi-omics data for discovery and functional studies of complex human diseases." *Advances in genetics*. Vol. 93. Academic Press, 2016. 147-190.
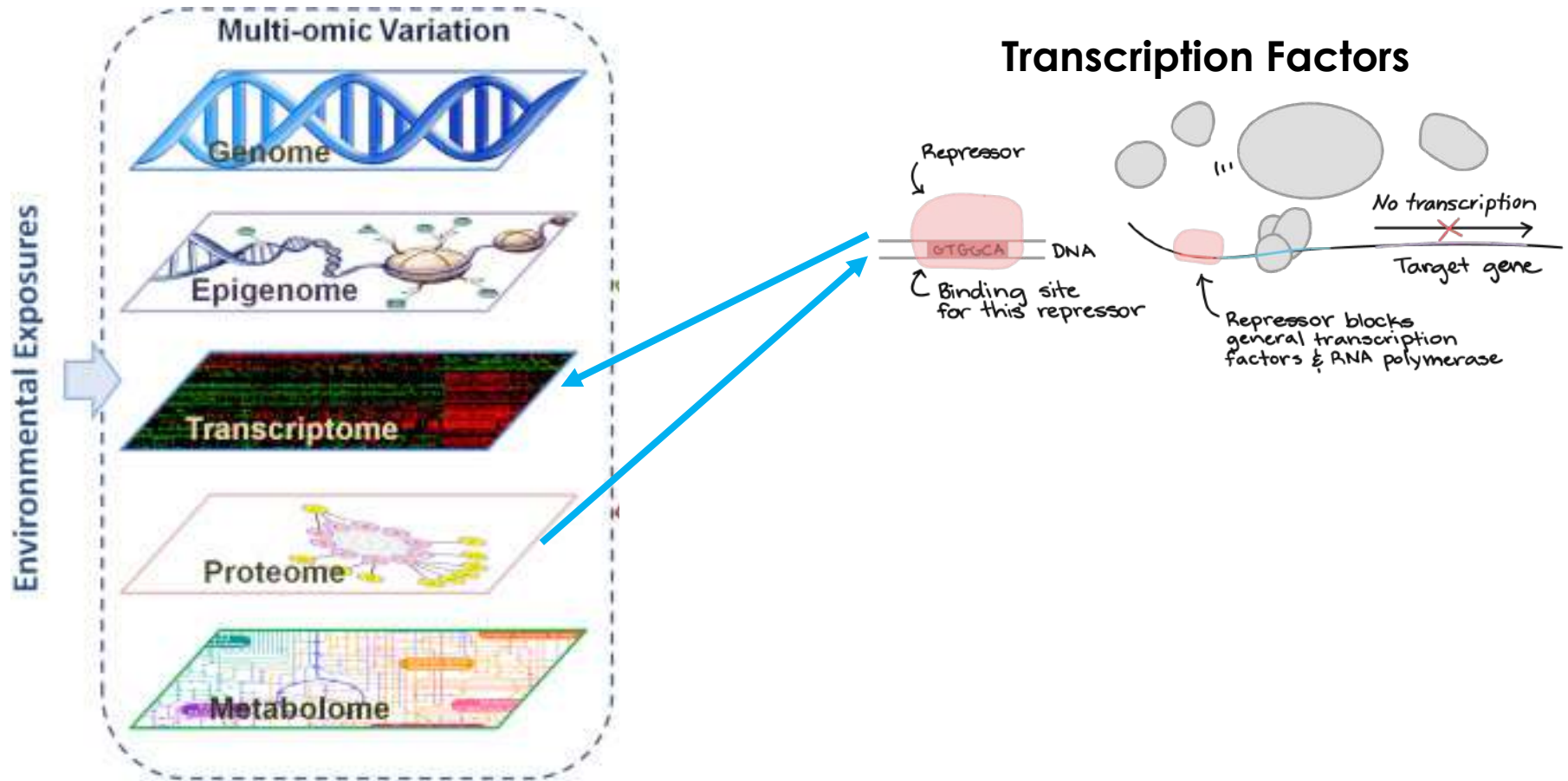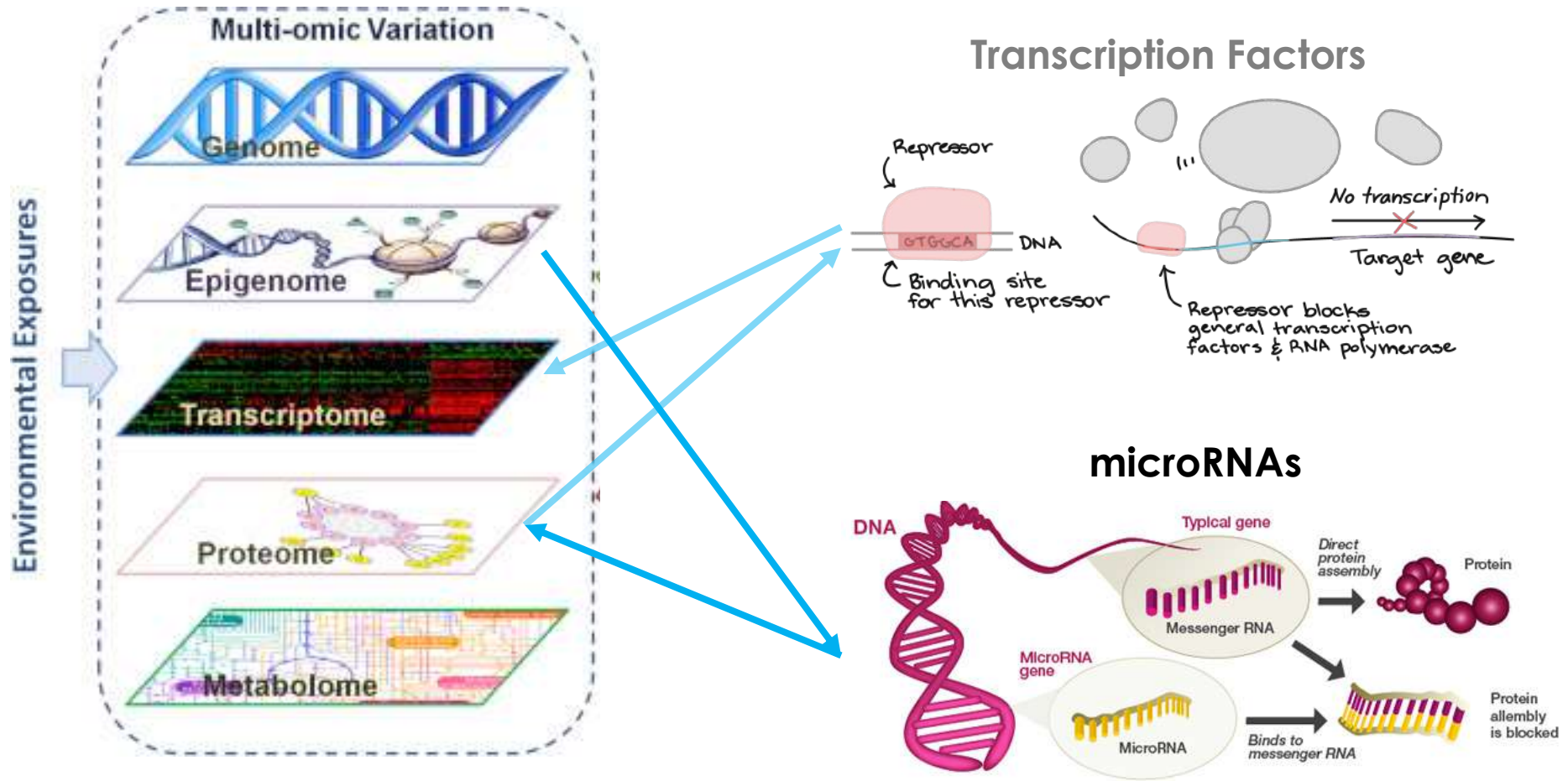
# Multi-omics data are interconnected

The Cancer Genome Atlas Research Network, Weinstein, J.N., Collisson, E.A., Mills, G.B., Shaw, K.M., Ozenberger, B.A., Ellrott, K., Shmulevich, I., Sander, C., and Stuart, J.M. (2013) The Cancer Genome Atlas Pan-Cancer analysis project. Nat Genet. doi:10.1038/ng.2764

Sun, Yan V., and Yi-Juan Hu. "Integrative analysis of multi-omics data for discovery and functional studies of complex human diseases." *Advances in genetics*. Vol. 93. Academic Press, 2016. 147-190.
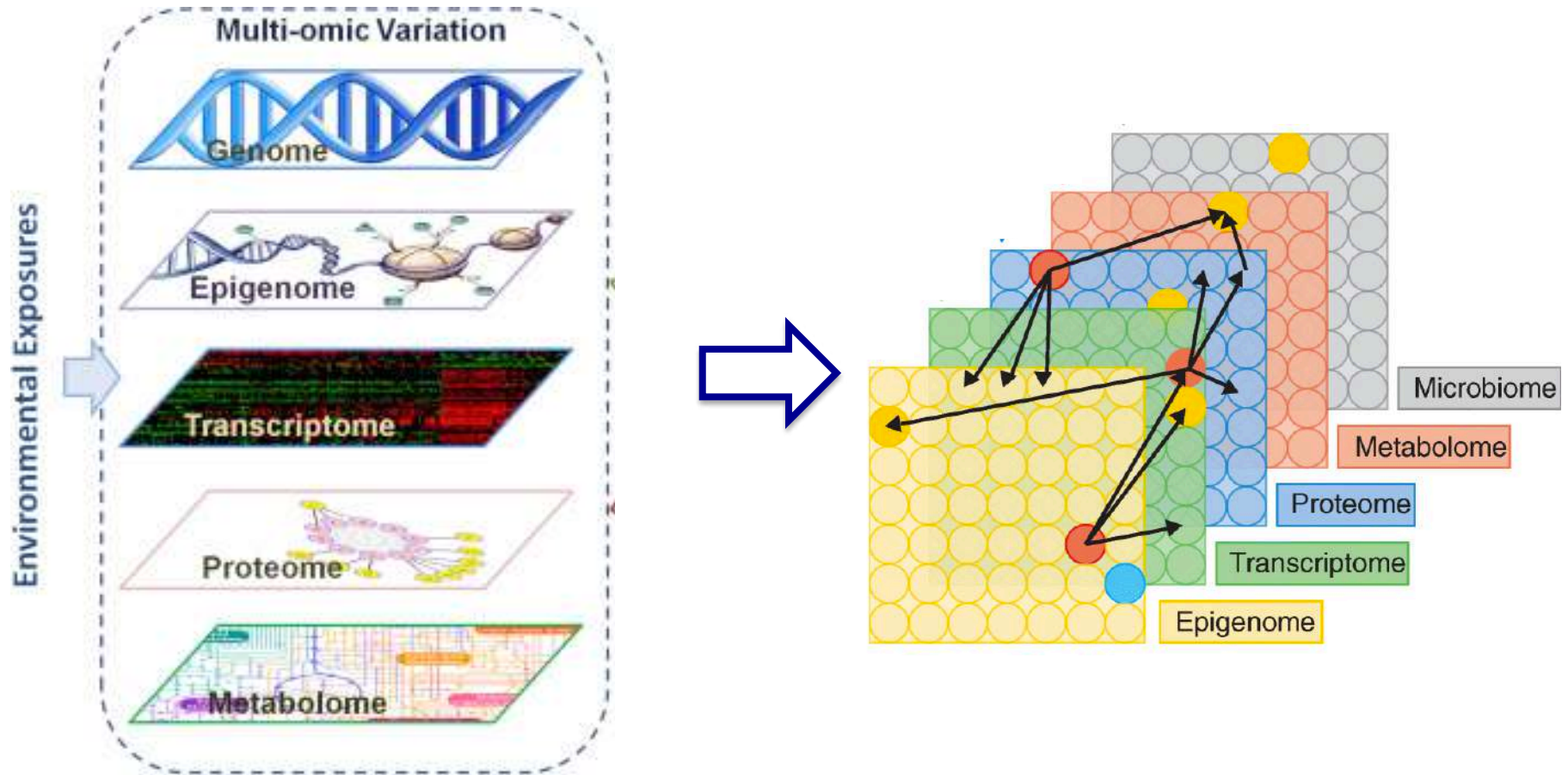
# Multi-omics data are interconnected

The Cancer Genome Atlas Research Network, Weinstein, J.N., Collisson, E.A., Mills, G.B., Shaw, K.M., Ozenberger, B.A., Ellrott, K., Shmulevich, I., Sander, C., and Stuart, J.M. (2013) The Cancer Genome Atlas Pan-Cancer analysis project. Nat Genet. doi:10.1038/ng.2764
Sun, Yan V., and Yi-Juan Hu. "Integrative analysis of multi-omics data for discovery and functional studies of complex human diseases." *Advances in genetics*. Vol. 93. Academic Press, 2016. 147-190.

# Multi-omics data are interconnected

The Cancer Genome Atlas Research Network, Weinstein, J.N., Collisson, E.A., Mills, G.B., Shaw, K.M., Ozenberger, B.A., Ellrott, K., Shmulevich, I., Sander, C., and Stuart, J.M. (2013) The Cancer Genome Atlas Pan-Cancer analysis project. Nat Genet. doi:10.1038/ng.2764

Sun, Yan V., and Yi-Juan Hu. "Integrative analysis of multi-omics data for discovery and functional studies of complex human diseases." *Advances in genetics*. Vol. 93. Academic Press, 2016. 147-190.

# Multi-omics data are interconnected



# The joint analysis of multiple omics is required

The Cancer Genome Atlas Research Network, Weinstein, J.N., Collisson, E.A., Mills, G.B., Shaw, K.M., Ozenberger, B.A., Ellrott, K., Shmulevich, I., Sander, C., and Stuart, J.M. (2013) The Cancer Genome Atlas Pan-Cancer analysis project. Nat Genet. doi:10.1038/ng.2764
Sun, Yan V., and Yi-Juan Hu. "Integrative analysis of multi-omics data for discovery and functional studies of complex human diseases." *Advances in genetics*. Vol. 93. Academic Press, 2016. 147-190.

# Challenges of multi-omics integration

High-dimensionality -> Big-data

Heterogeneous variables

Different ranges of variation

Technical noise different for each omics

# More omics is better, but how many more?

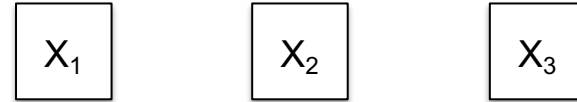# Is it always good to consider ALL the available omics?
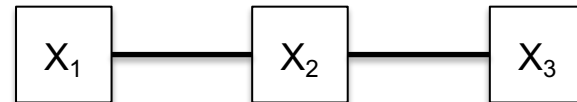
# Choosing which omics to integrate
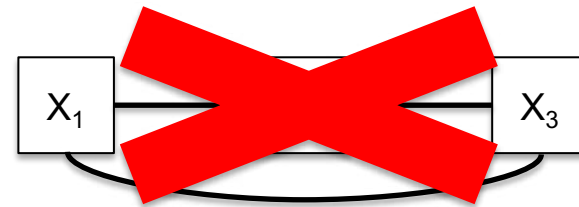
**Aim:** predicting drug response

**Available input data:**

- Mutations

- Copy Number Alterations (CNA)

- Methylation

- Gene expression
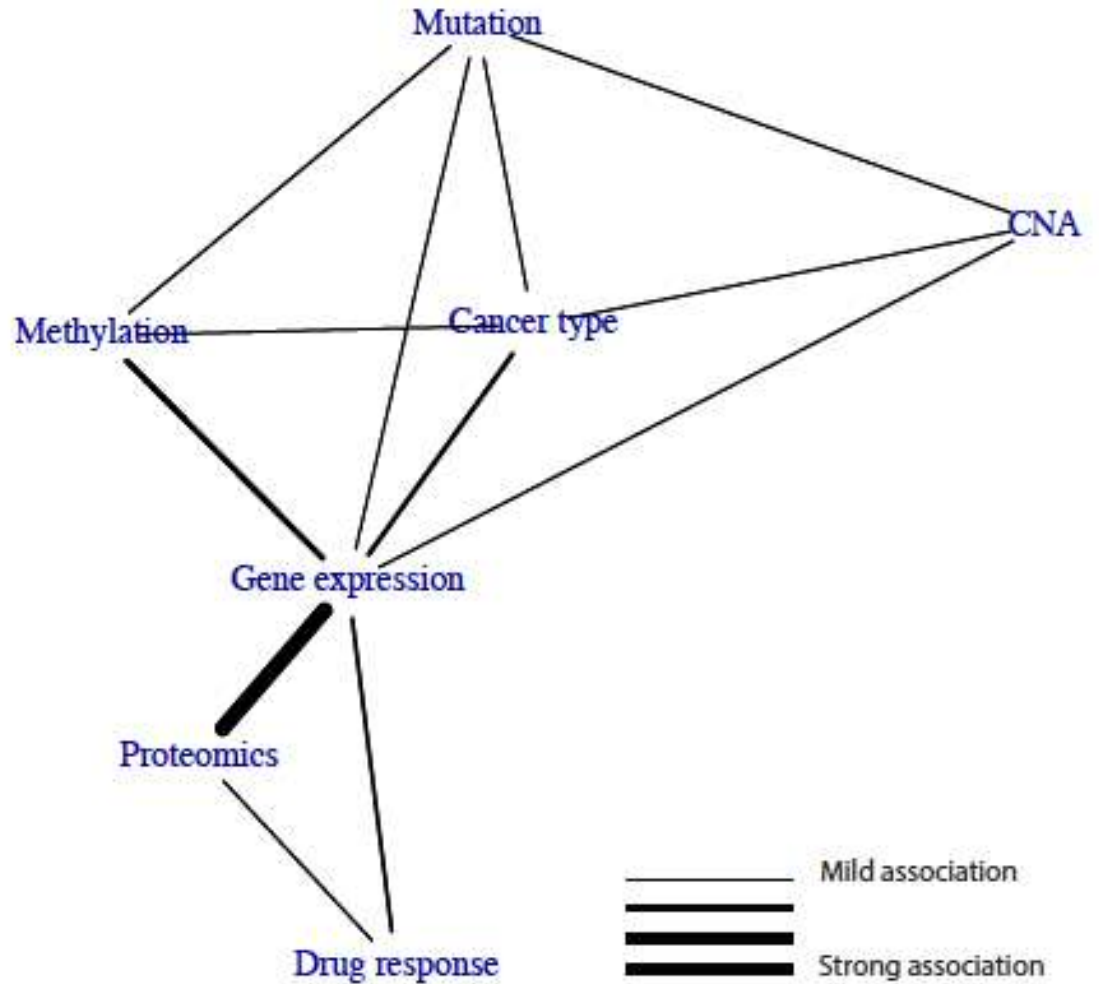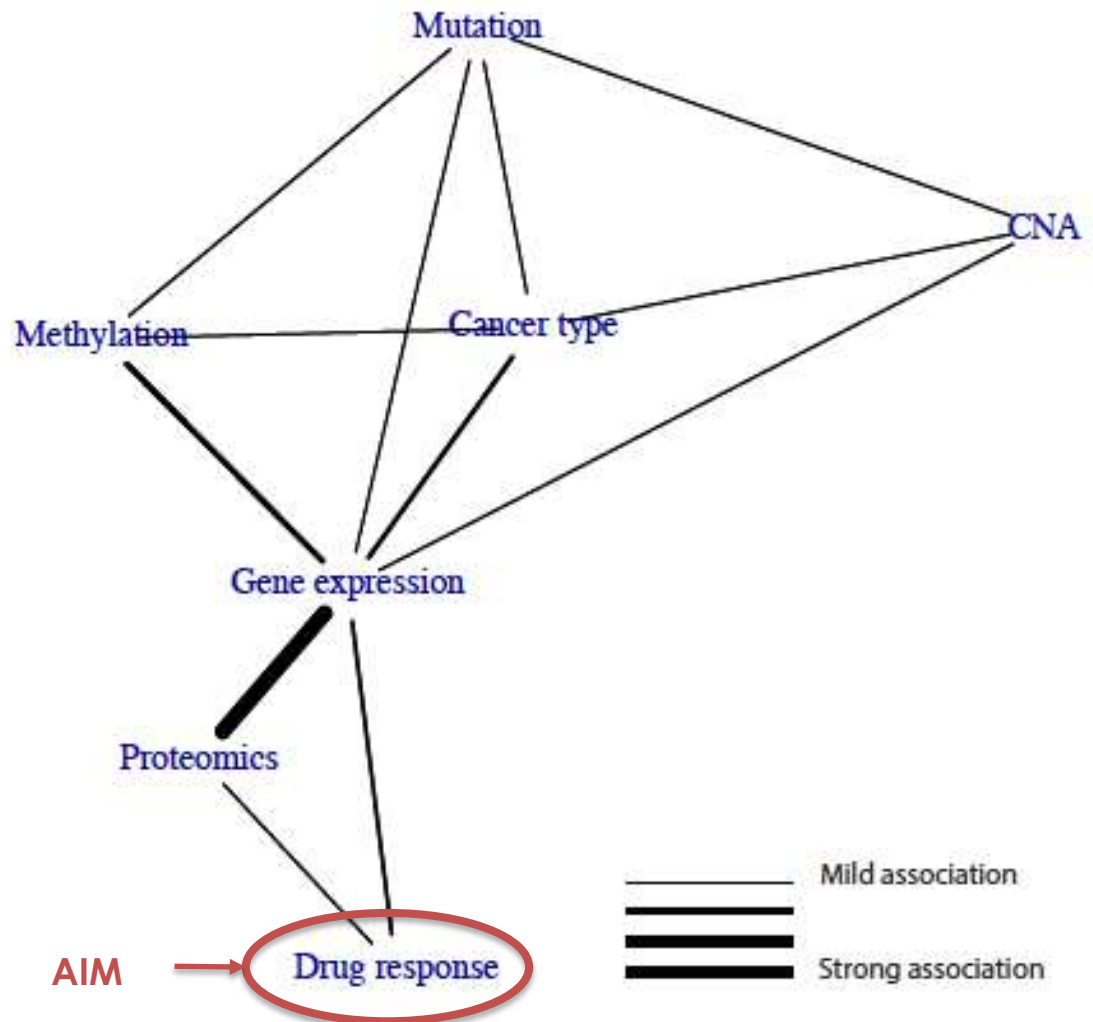
- Proteomics

- Cancer types

- Drug response

ABEN, Nanne, et al. iTOP: inferring the topology of omics data. Bioinformatics, 2018, 34.17: i988-i996.

# Choosing which omics to integrate

**Aim:** predicting drug response

**Available input data:**

- Mutations
- Copy Number Alterations (CNA)
- Methylation
- Gene expression
- Proteomics
- Cancer types
- Drug response

$X_1$  $X_2$  $X_3$

**Using correlation:**



ABEN, Nanne, et al. iTOP: inferring the topology of omics data. Bioinformatics, 2018, 34.17: i988-i996.

# Choosing which omics to integrate

**Aim:** predicting drug response

**Available input data:**

- Mutations

- Copy Number Alterations (CNA)

- Methylation

- Gene expression

- Proteomics

- Cancer types

- Drug response

$X_1$ $X_2$ $X_3$

**Using partial correlation (iTOP):**

e.g. $cor(X_1, X_3 \mid X_2) \approx 0$



$X_1$ — $X_2$ — $X_3$

ABEN, Nanne, et al. iTOP: inferring the topology of omics data. Bioinformatics, 2018, 34.17: i988-i996.

# Choosing which omics to integrate



ABEN, Nanne, et al. iTOP: inferring the topology of omics data. Bioinformatics, 2018, 34.17: i988-i996.

# Choosing which omics to integrate

ABEN, Nanne, et al. iTOP: inferring the topology of omics data. Bioinformatics, 2018, 34.17: i988-i996.

# Choosing which omics to integrate

ABEN, Nanne, et al. iTOP: inferring the topology of omics data. Bioinformatics, 2018, 34.17: i988-i996.

How the omics should be combined?

# Integrating multi-omics data

**Approach "Genome First"**

Priority given to genome

Other omics are only used for interpretation

Hasin, Yehudit, Marcus Seldin, and Aldons Lusis. "Multi-omics approaches to disease." Genome biology 18.1 (2017): 83.

# Integrating multi-omics data



**Dataset 1**   **Dataset 2**

Machine learning algorithm designed for a single dataset

**How do I integrate them ???????????**

Zitnik, Marinka, et al. "Machine learning for integrating data in biology and medicine: Principles, practice, and opportunities." Information Fusion 50 (2019): 71-91.

# Integrating multi-omics data



**Late integration**
output averaging, ensembles

Dataset 1   Dataset 2

Outputs, predictions

machine learning model

Zitnik, Marinka, et al. "Machine learning for integrating data in biology and medicine: Principles, practice, and opportunities." Information Fusion 50 (2019): 71-91.

# Integrating multi-omics data

Zitnik, Marinka, et al. "Machine learning for integrating data in biology and medicine:
Principles, practice, and opportunities." Information Fusion 50 (2019): 71-91.

# Integrating multi-omics data



**Late integration**
output averaging, ensembles

**Early integration**
projection, concatenation
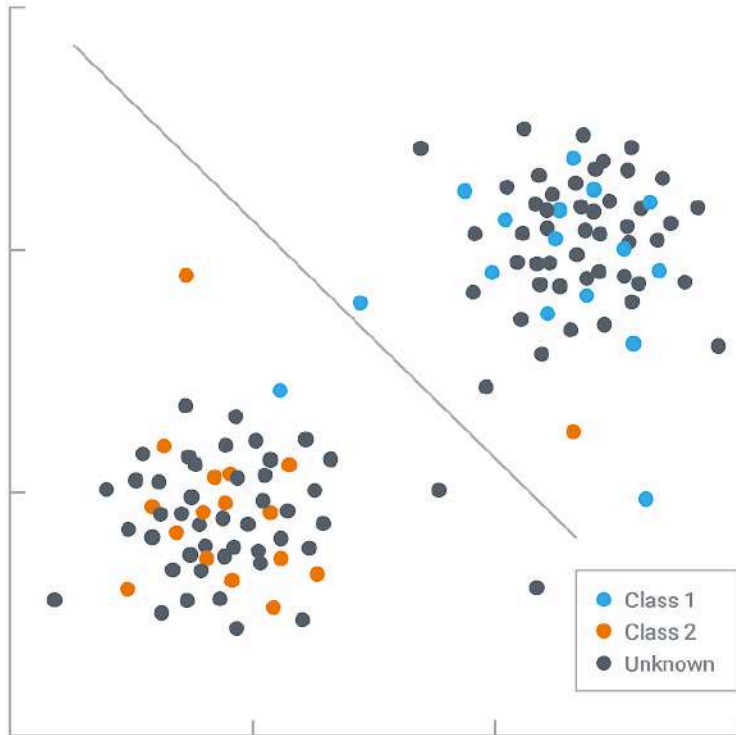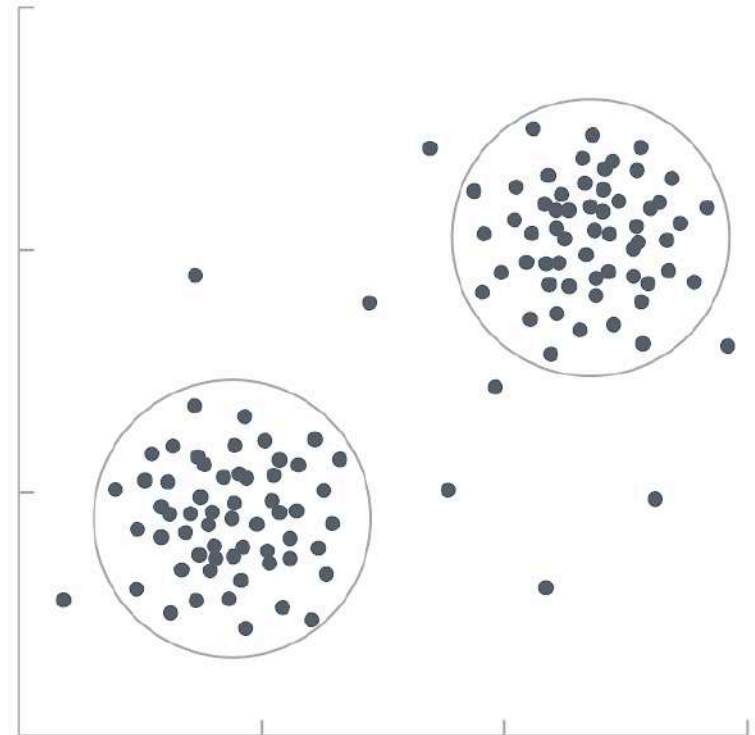
**Intermediate integration**
multi-view, multi-modal

Dataset 1  Dataset 2

Combined dataset

Outputs, predictions

machine learning model

Zitnik, Marinka, et al. "Machine learning for integrating data in biology and medicine: Principles, practice, and opportunities." Information Fusion 50 (2019): 71-91.

Main categories of existing multi-omics integrative approaches

# Main categories of integrative approaches

## Supervised methods



Class 1
Class 2
Unknown

## Unsupervised methods

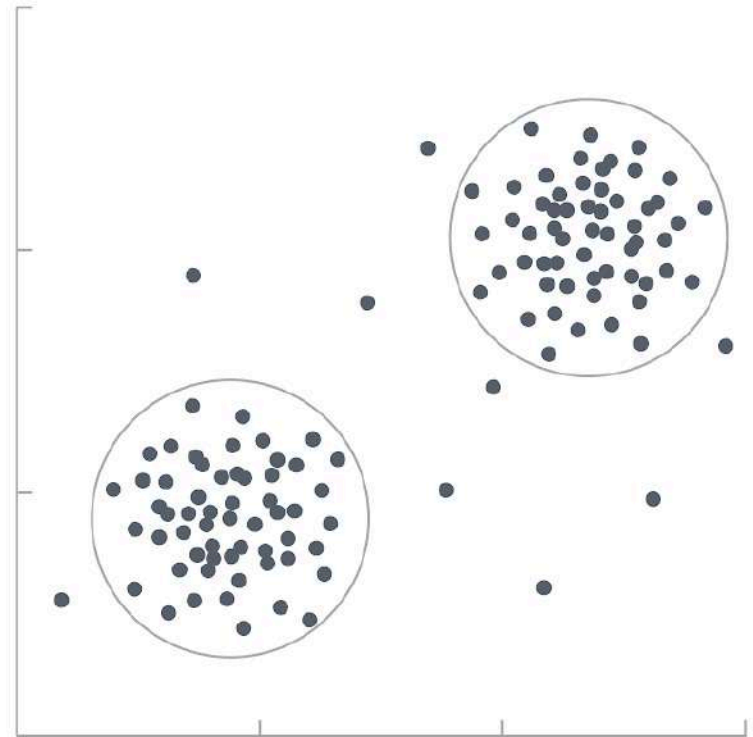# Main categories of integrative approaches

## Supervised methods



- They require 2 datasets in input: training and test datasets

- Labels must be avilable for the training dataset

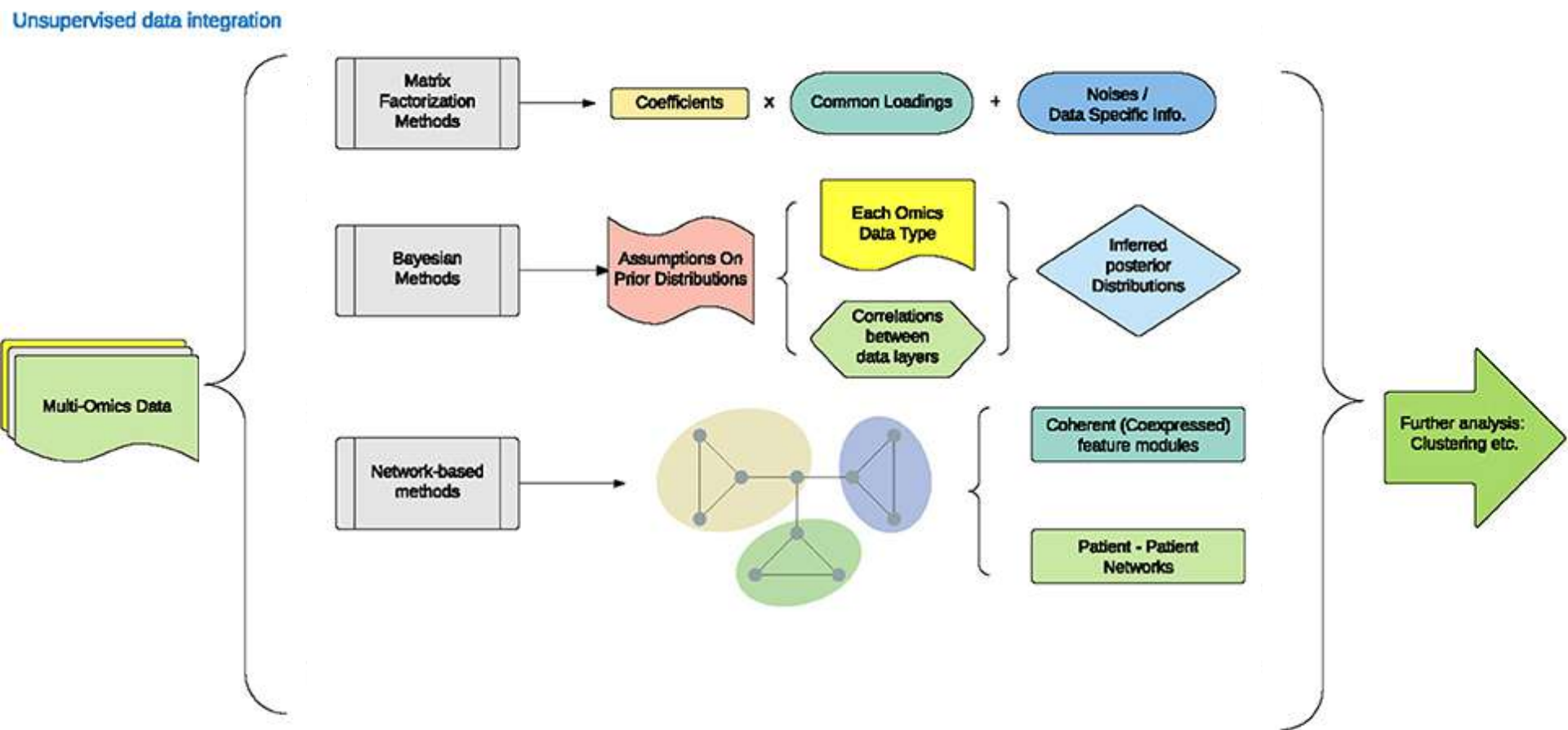- This information is used to infer labels on the test dataset

# Main categories of integrative approaches

## Unsupervised methods

- The methodology is directly applied to one dataset

- They infer information from the structure of the data without any label information
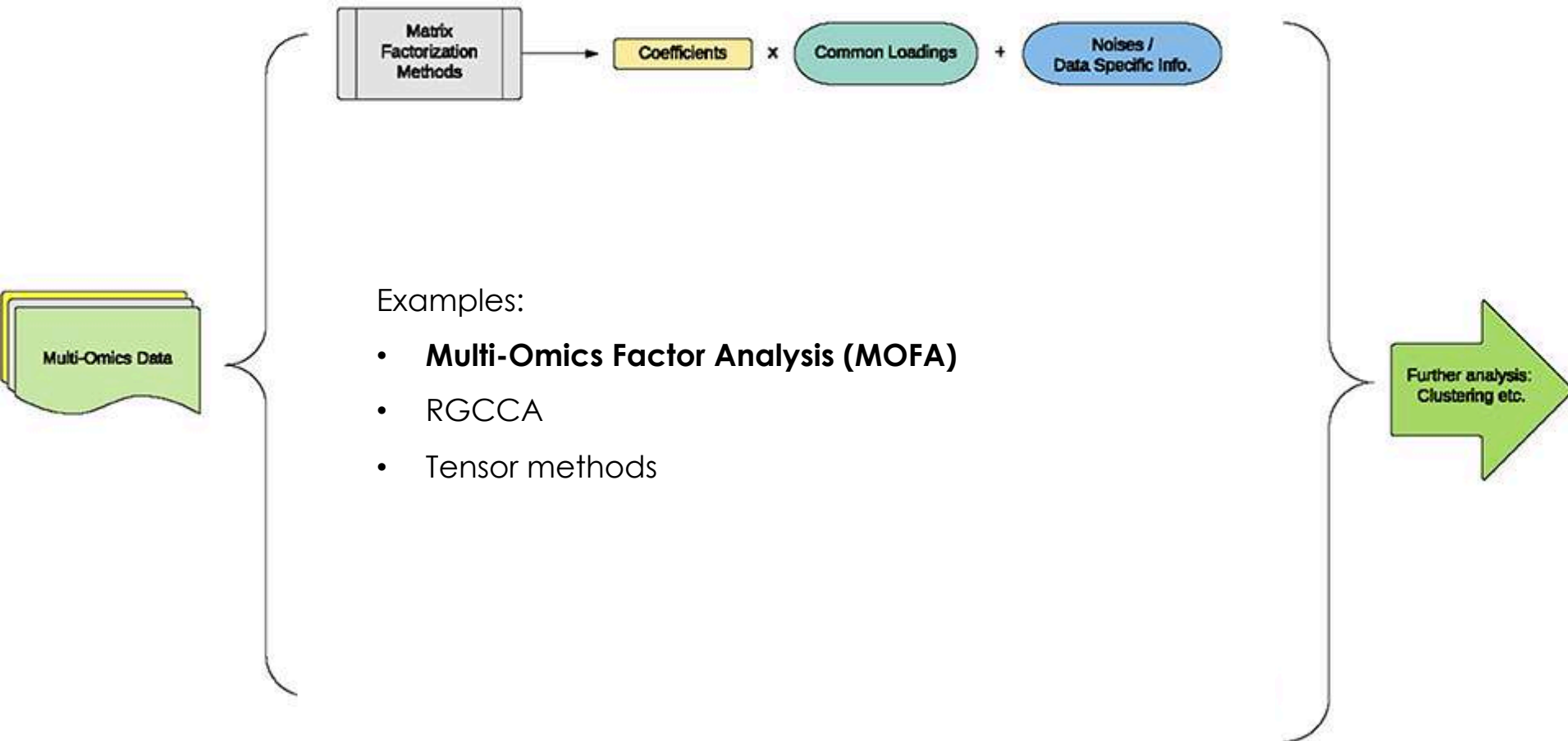
# Unsupervised integrative approaches

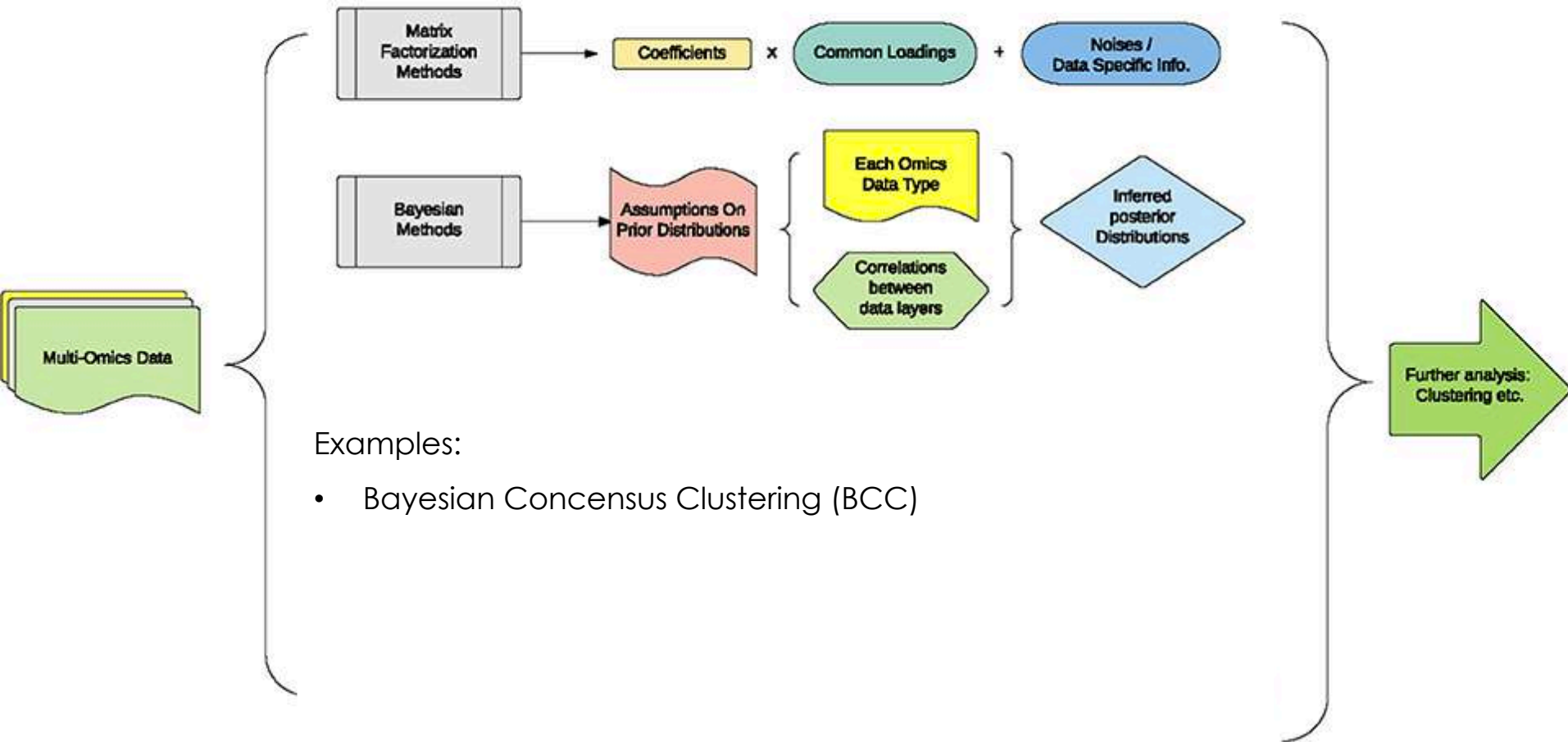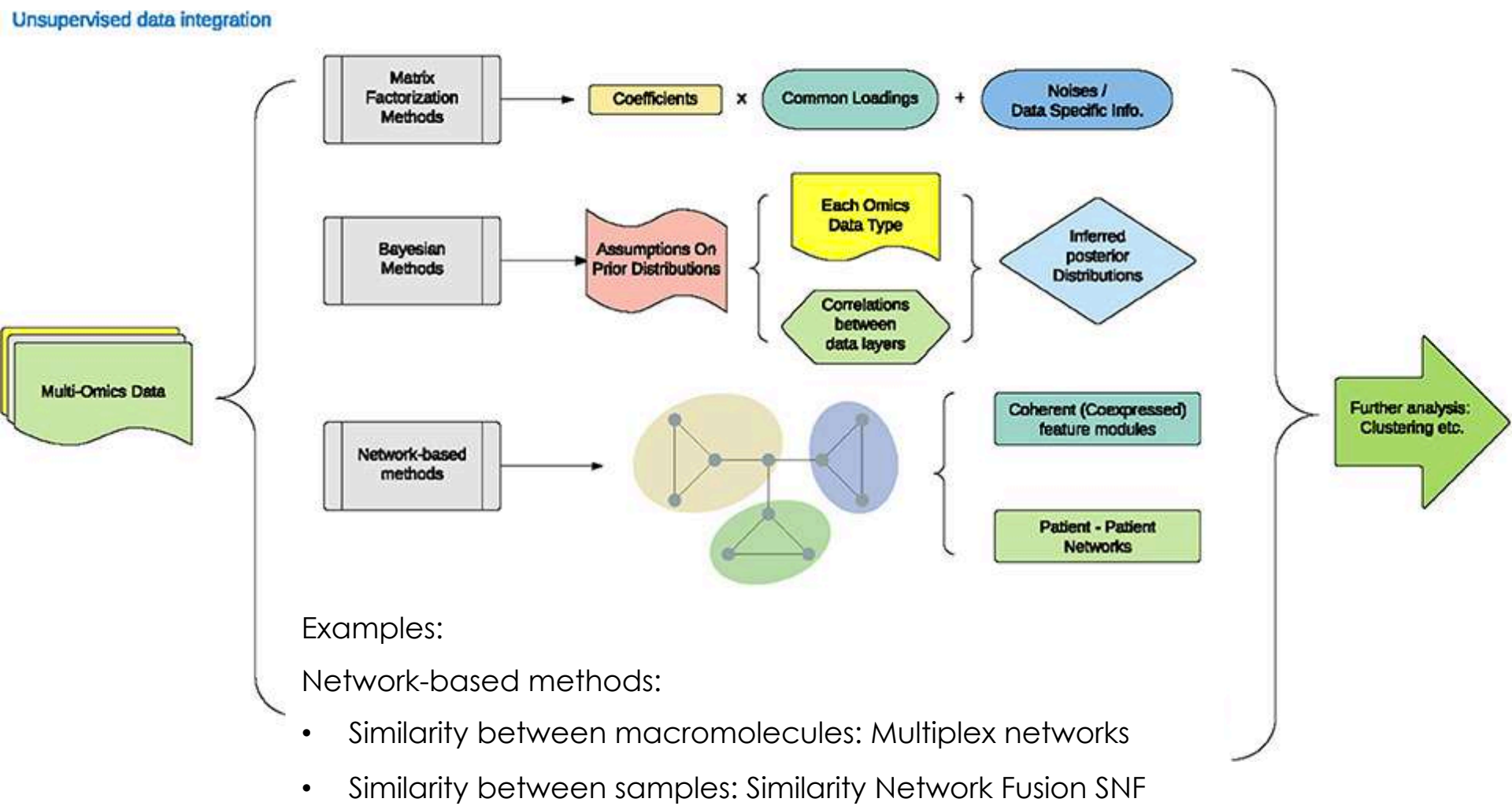# Unsupervised integrative approaches

Unsupervised data integration



Examples:

- **Multi-Omics Factor Analysis (MOFA)**

- RGCCA

- Tensor methods

# Unsupervised integrative approaches



Examples:

- Bayesian Concensus Clustering (BCC)

# Unsupervised integrative approaches



Examples:

Network-based methods:

- Similarity between macromolecules: Multiplex networks
- Similarity between samples: Similarity Network Fusion SNF

# Cancer insights from data integration methods

# Cancer subtyping



| CMS1 (13%) | CMS2 (35%) | CMS3 (11%) | CMS4 (20%) | Unclassified (21%) |
|---|---|---|---|---|
| • Right colon, female<br>• MSI, *BRAF* mut, hypermutated<br>• Immune activation<br>• Worse survival after relapse | • Left colon<br>• MSS, CIN, *BRAF* wt, *TP53* mut<br>• Epithelial, WNT/Myc pathway activation<br>• Better survival after relapse | • *KRAS* mut<br>• Epithelial, IGFBP2 overexpression | • Mesenchymal, TGFβ pathway activation, NOTCH3 overexpression<br>• Worse relapse free survival and overall survival | • Immune and stroma infiltration<br>• Variable epithelial - mesenchymal activation |

CMS1: C2, Subtype 1.2, A-type, CCS2, C, Inflammatory

CMS2: C1-C5-C6, B-type, Subtype 2.2, CCS1, B, Enterocyte-TA

CMS3: C3, Subtype 2.1, Globet-like, A

CMS4: C4, C-type, Subtype 1.1-1.3, CCS3, D-E, Stem-like

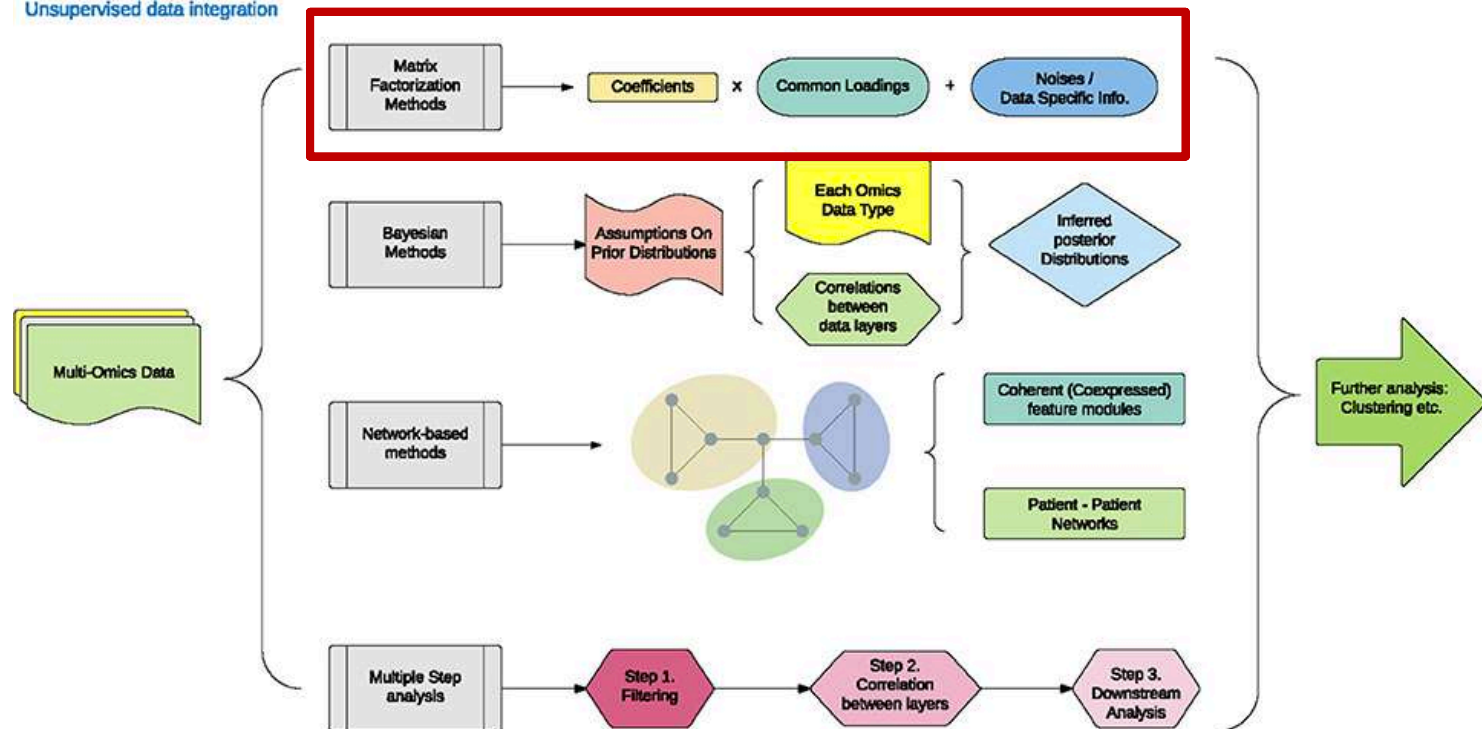Santos, Cristina, et al. "Intrinsic cancer subtypes-next steps into personalized medicine." Cellular oncology 38.1 (2015): 3-16.

# Cancer subtyping

This problem is generally approached with unsupervised approaches.

# Gene modules identification

Drug rensponding

Drug resistent



Which are the molecular mechanisms that make these two groups of patients having a different behaviour?

Can we identify a driver that can alter the behaviour of a set of patients?

# Gene modules identification

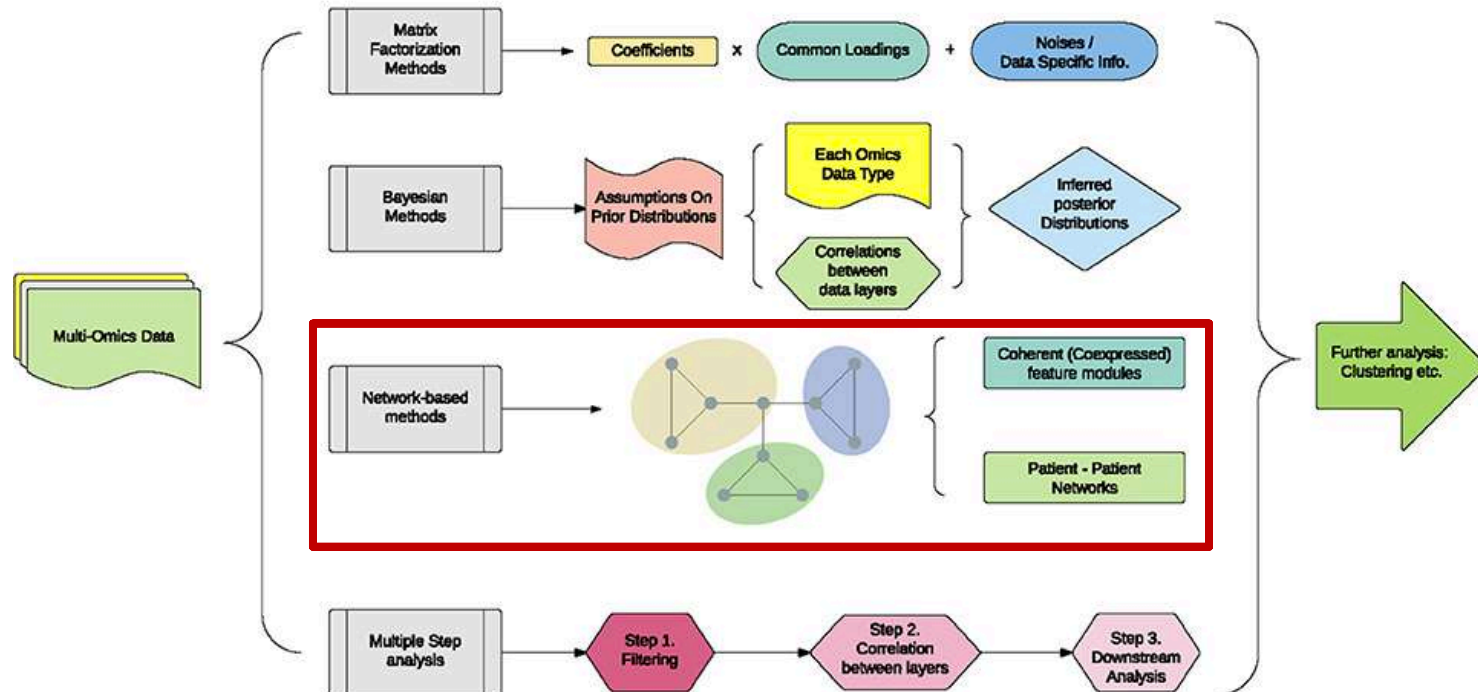This problem is generally approached with unsupervised approaches.
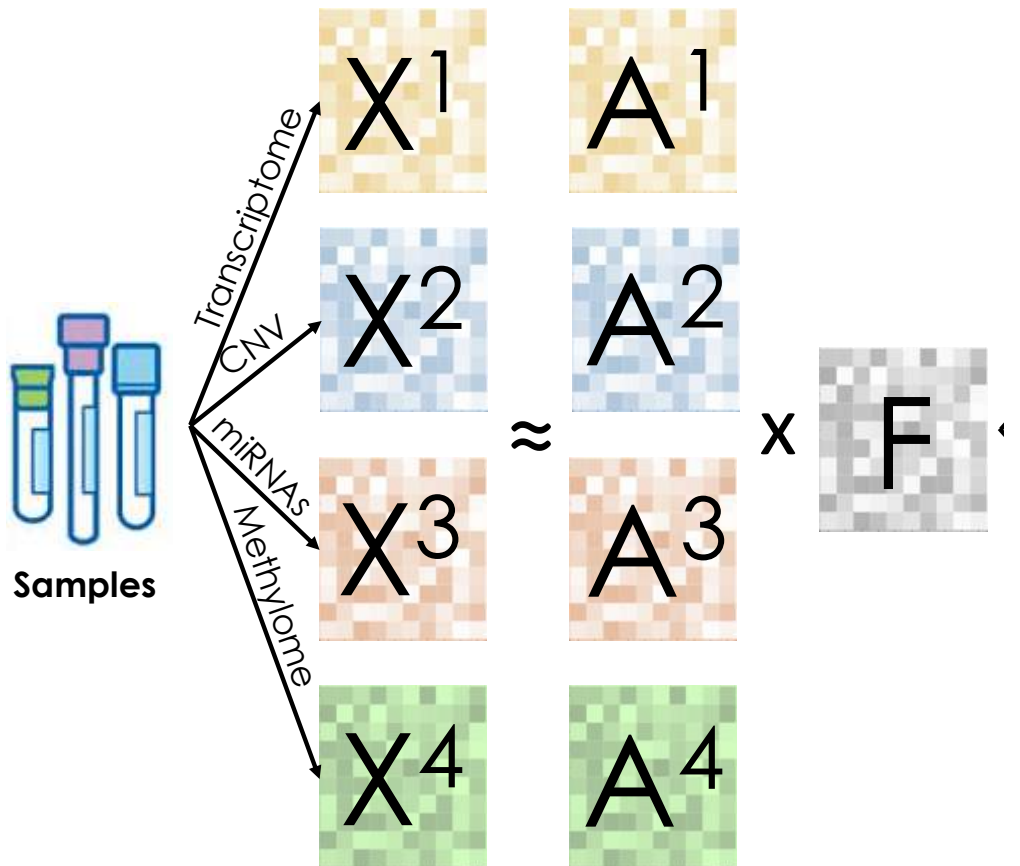
Drug rensponding

Drug resistent

# Matrix Factorization

# Joint Dimensionality Reduction (jDR)



$$X^m = A^m F + e^m$$
$$e^m <<<<$$
$$m = 1, \ldots, P$$

Meng, Chen, et al. "Dimension reduction techniques for the integrative analysis of multi-omics data." *Briefings in bioinformatics* 17.4 (2016): 628-641.

# Joint Dimensionality Reduction (jDR)

**Multi-omics joint Dimensionality Reduction (jDR)**



**Samples**

Transcriptome
CNV
miRNAs
Methylome

$X^1$ $A^1$

$X^2$ $A^2$

$\approx$ $X^3$ $A^3$ x $F$

$X^4$ $A^4$

**Sample clustering**

Gene 1
Gene 13 Gene 3 Gene 20
Gene 2 Gene 10
Gene 11

**Pathways/processes/markers/
molecular mechanisms**

$$X^m = A^m F + e^m$$
$$e^m \lll$$
$$m = 1, \ldots, P$$

Meng, Chen, et al. "Dimension reduction techniques for the integrative analysis of multi-omics data." *Briefings in bioinformatics* 17.4 (2016): 628-641.

# Multi-omics Factor Analysis (MOFA)



$$Y^m = ZW^m + e^m$$
$$m = 1, \ldots, M$$

Argelaguet, Ricard, et al. "Multi-Omics Factor Analysis—a framework for unsupervised integration of multi-omics data sets." Molecular systems biology 14.6 (2018): e8124.

# MOFA advantage: interpretability of factors

Argelaguet, Ricard, et al. "Multi-Omics Factor Analysis—a framework for unsupervised integration of multi-omics data sets." Molecular systems biology 14.6 (2018): e8124.

# Also single-cell multi-omics data can be integrated with matrix factorization

# Multi-omics single-cell



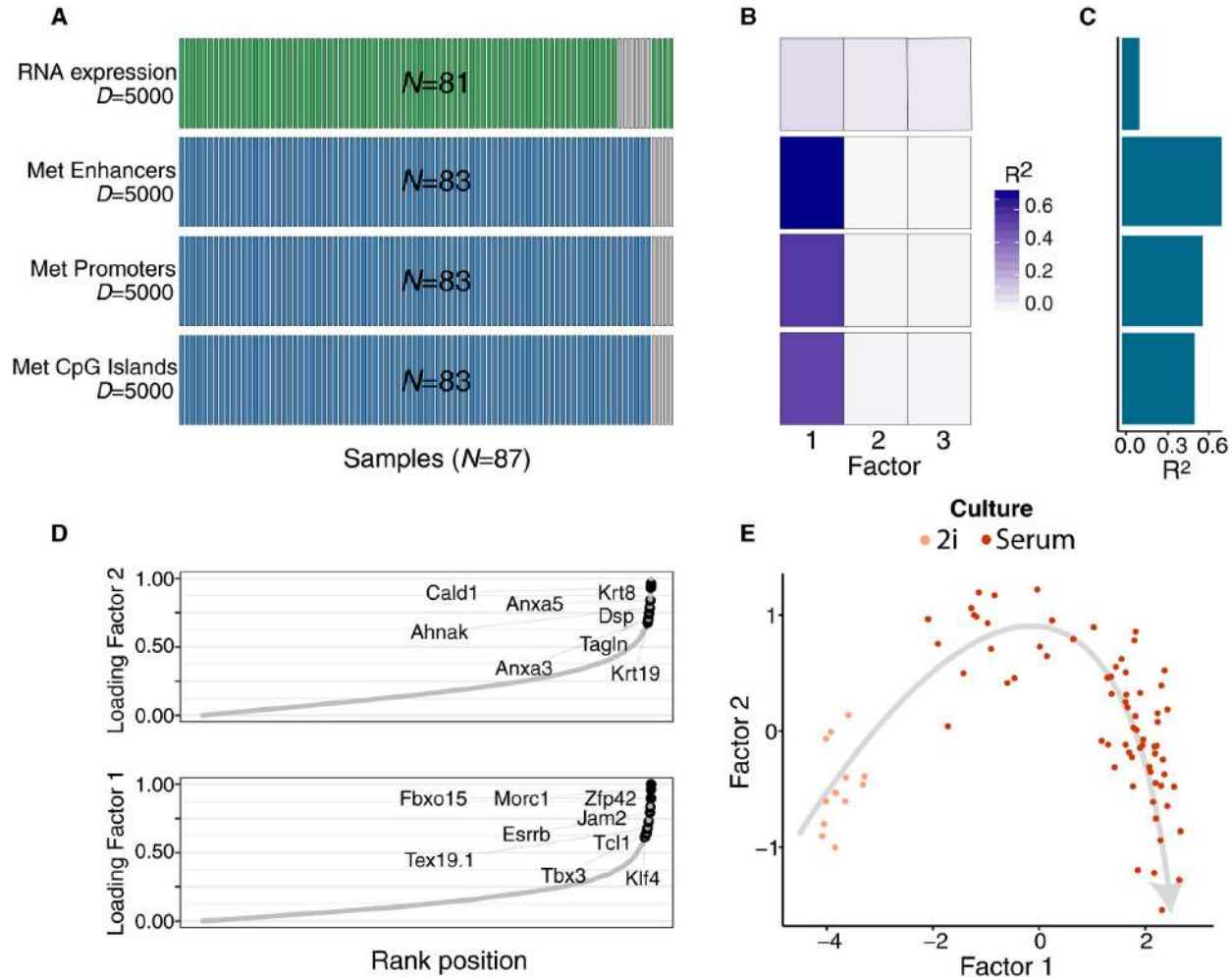Hu, Youjin, et al. " *Frontiers in cell and developmental biology* 6 (2018): 28.

# Example MOFA application single-cell multi-omics

Dataset: 87 mouse embryonic stem cells (mESCs) comprising:

- 16 cells cultured in "2i" media, which induces a naive pluripotency state

- 71 serum-grown cells, which commits cells to a primed pluripotency state poised for cellular differentiation.
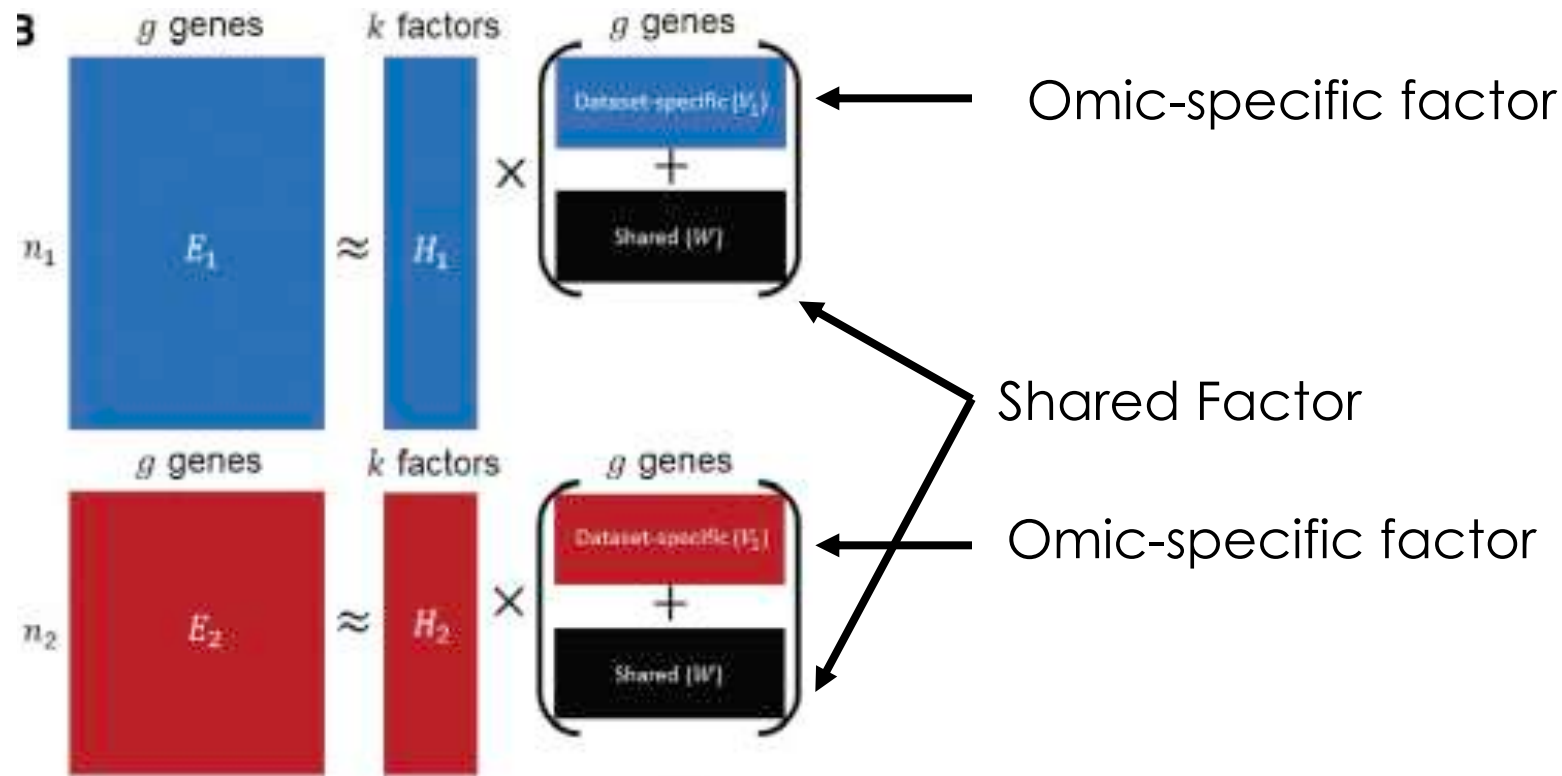
All cells were profiled using single-cell methylation and transcriptome sequencing

Argelaguet, Ricard, et al. "Multi-Omics Factor Analysis—a framework for unsupervised integration of multi-omics data sets." Molecular systems biology 14.6 (2018): e8124.

# Example MOFA application single-cell multi-omics

Argelaguet, Ricard, et al. "Multi-Omics Factor Analysis—a framework for unsupervised integration of multi-omics data sets." Molecular systems biology 14.6 (2018): e8124.

# Linked inference of genomic experimental relationships (LIGER)

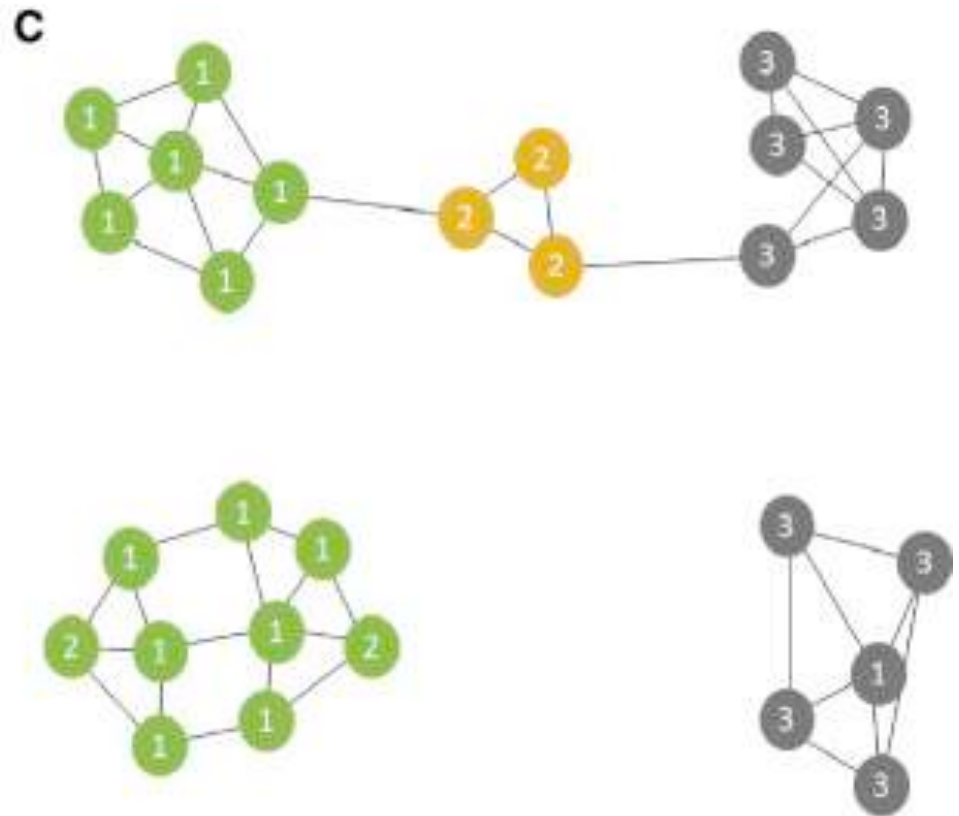**Integrative Non Negative Matrix Factorization (iNMF)**



$$E_i = H_i \, V_i + H_i \, W$$

Welch, Joshua D., et alCell 177.7 (2019): 1873-1887.

# LIGER: multi-omics clustering

**Integrative Non Negative Matrix Factorization (iNMF)**

**kNN graphs to derive clusters from factors**



$$E_i = H_i V_i + H_i W$$

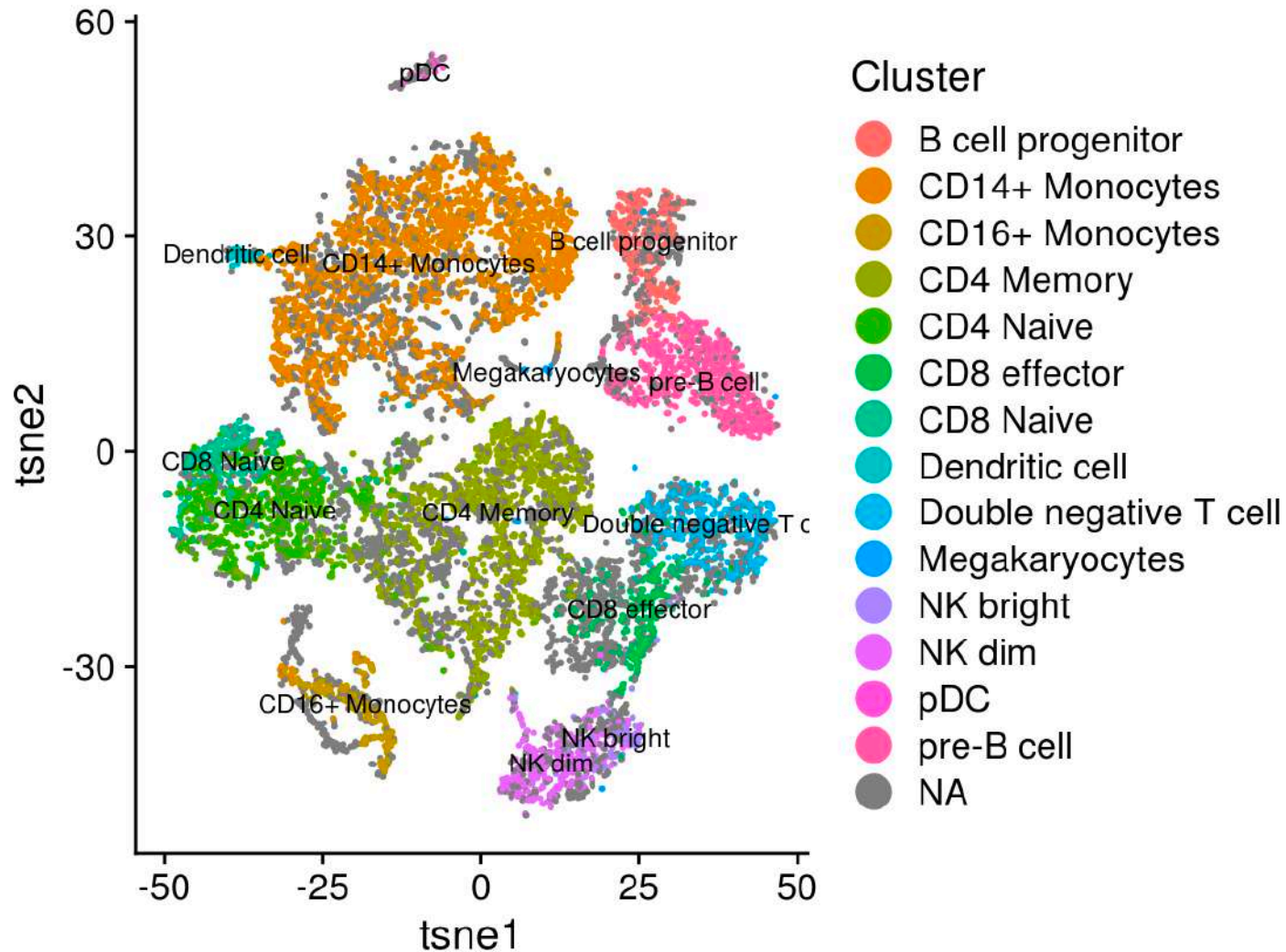Welch, Joshua D., et alCell 177.7 (2019): 1873-1887.

# LIGER: peripheral blood mononuclear cell (PBMC)

scRNAseq and scATACseq data from approx. 10k cells PBMCs

We want to identify subtypes of cells based on the joint analysis of the two data types

Welch, Joshua D., et alCell 177.7 (2019): 1873-1887.

# LIGER: peripheral blood mononuclear cell (PBMC)



Welch, Joshua D., et alCell 177.7 (2019): 1873-1887.
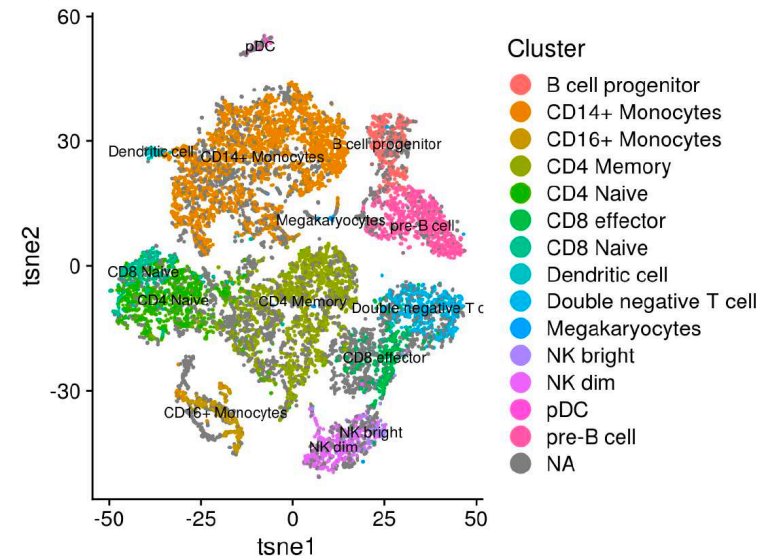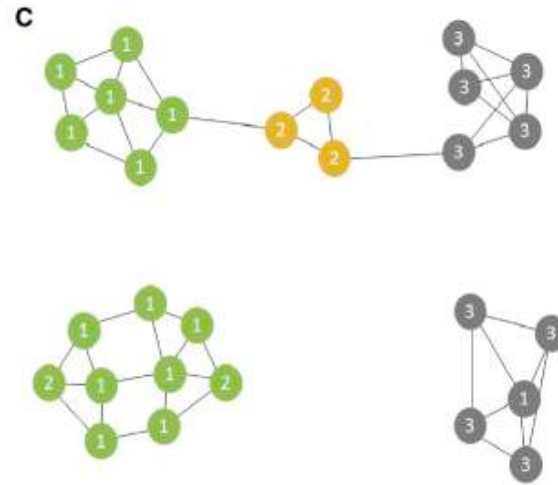
# Pay attention this is not a TSNE plot of scRNAseq data



**Integrative Non Negative Matrix Factorization (iNMF)**

**kNN graphs to derive clusters from factors**

$$E_i = H_i V_i + H_i W$$

Welch, Joshua D., et al Cell 177.7 (2019): 1873-1887.

# LIGER: peripheral blood mononuclear cell (PBMC)



Welch, Joshua D., et alCell 177.7 (2019): 1873-1887.