

Data integration in cancer research

An overview of the existing approaches

Laura Cantini

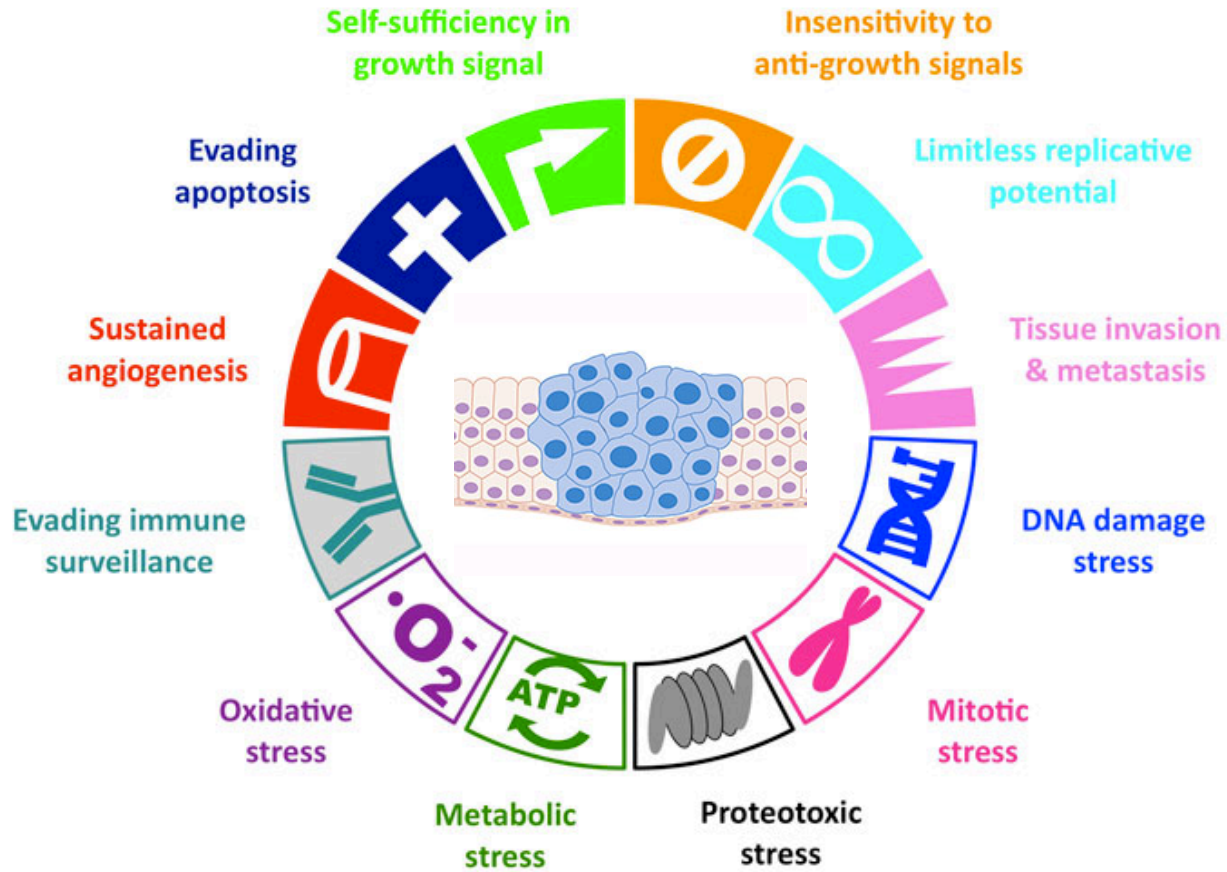
Computational Systems Biology Team
IBENS, Paris

26-02-2019

DU-Bii Integrative bioinformatics

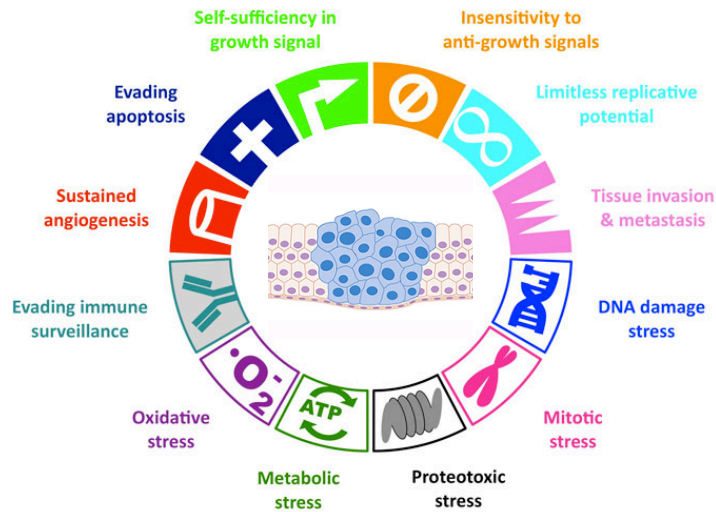


Cancer Hallmarks



Hanahan, Douglas, and Robert A. Weinberg. "Hallmarks of cancer: the next generation." *cell* 144.5 (2011): 646-674.

Cancer Hallmarks



2012→ 2030

**WORLDWIDE CANCER CASES
ARE PROJECTED TO INCREASE BY**

↑ 50%

FROM 14 million TO 21 million

**WORLDWIDE CANCER DEATHS
ARE PROJECTED TO INCREASE BY**

↑ 63%

FROM 8 million TO 13 million

Source: American Cancer Society: Global Cancer Facts & Figures, Second Edition

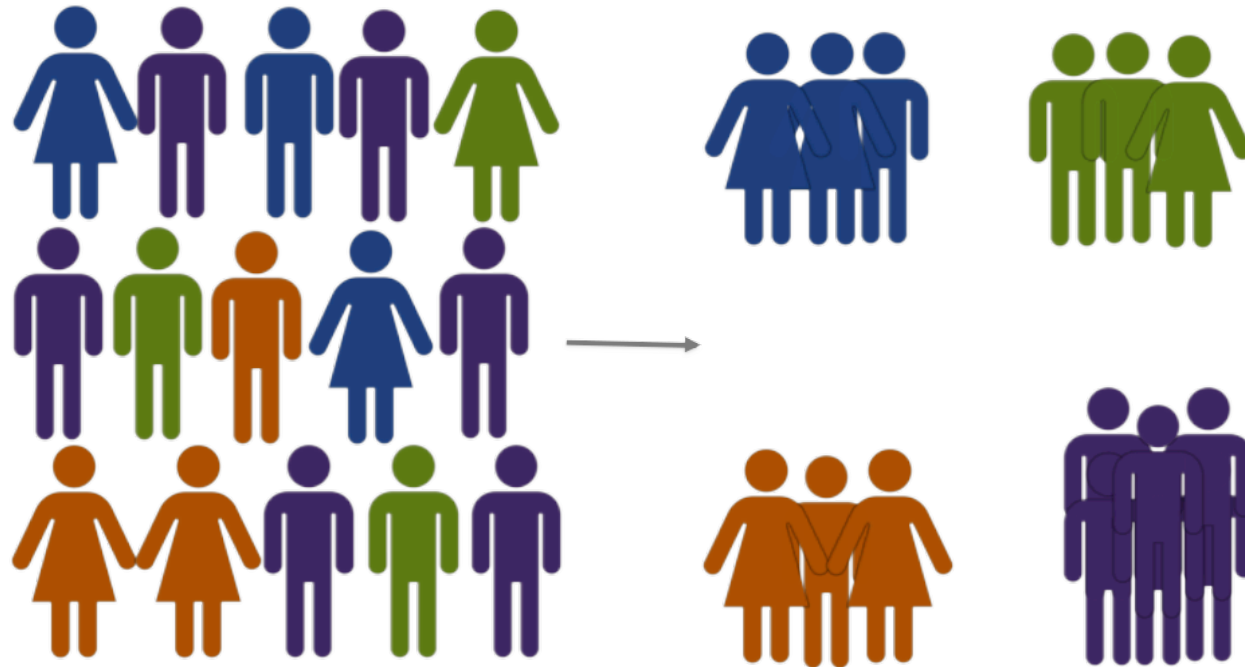
NIH
<https://www.cancer.gov/>

Hanahan, Douglas, and Robert A. Weinberg. "Hallmarks of cancer: the next generation." cell 144.5 (2011): 646-674.

Cancer precision medicine

Patients suffering from the same cancer can present:

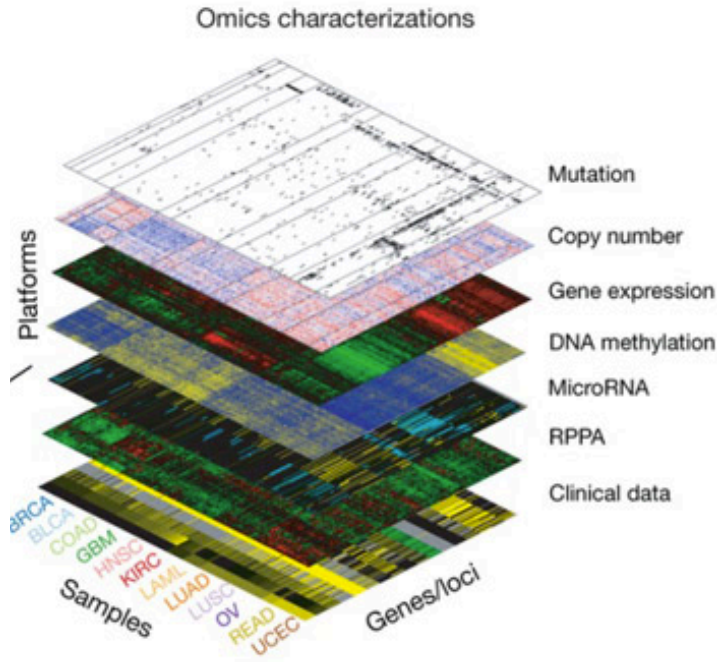
- Different prognosis
- Different response to the same treatment



→ Precision medicine is needed

Cancer « omics » data

The Cancer Genome Atlas (TCGA)

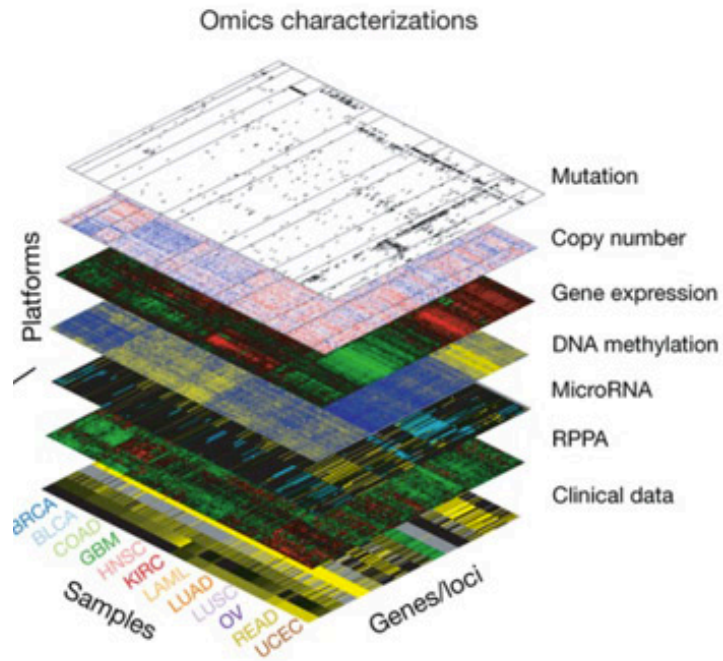


From single omics to **Multi-omics**

TCGA is the largest collection of multi-omics data:

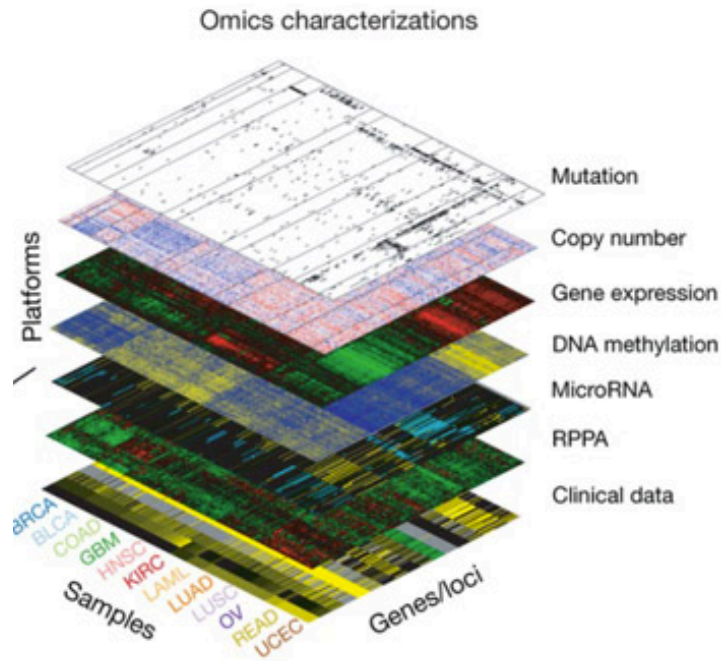
- 10.000 cancer patients
- 33 cancer types
- 6 omics, plus clinical data

Multi-omics integration in cancer

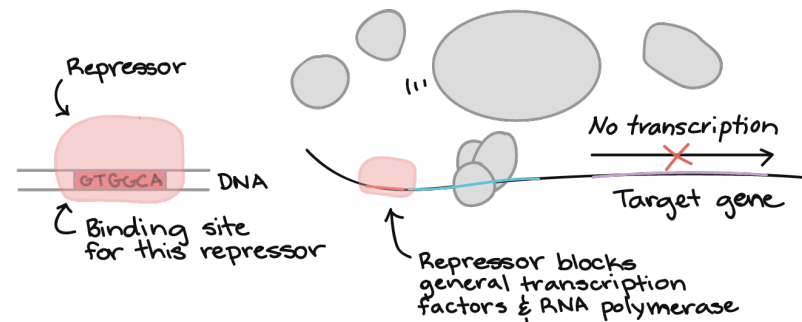


The Cancer Genome Atlas Research Network, Weinstein, J.N., Collisson, E.A., Mills, G.B., Shaw, K.M., Ozenberger, B.A., Ellrott, K., Shmulevich, I., Sander, C., and Stuart, J.M. (2013) The Cancer Genome Atlas Pan-Cancer analysis project. Nat Genet. doi: 10.1038/ng.2764

Multi-omics integration in cancer

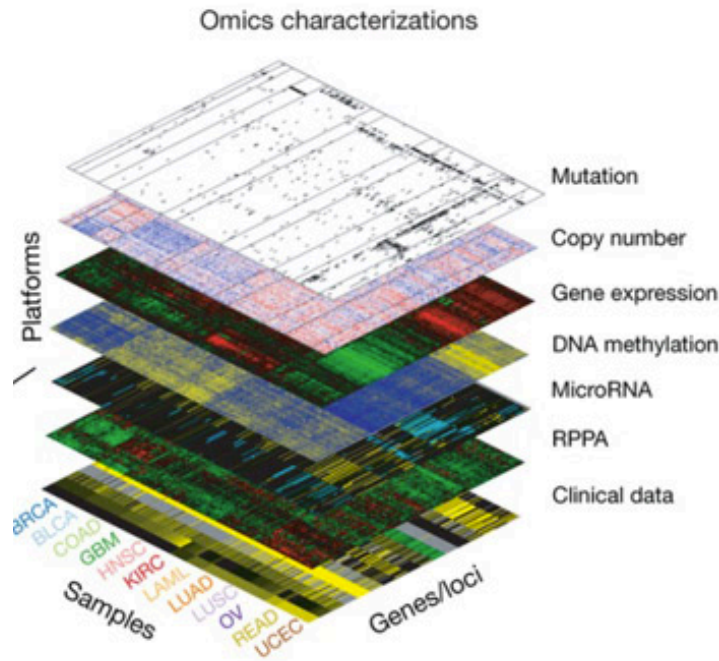


RPPA → Gene expression
(Transcription Factors)

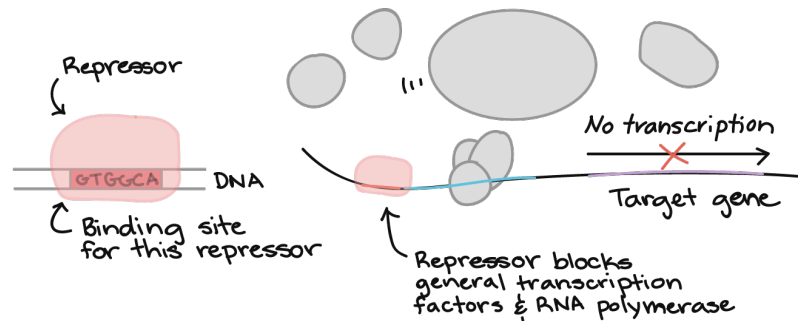


The Cancer Genome Atlas Research Network, Weinstein, J.N., Collisson, E.A., Mills, G.B., Shaw, K.M., Ozenberger, B.A., Ellrott, K., Shmulevich, I., Sander, C., and Stuart, J.M. (2013) The Cancer Genome Atlas Pan-Cancer analysis project. Nat Genet. doi: 10.1038/ng.2764

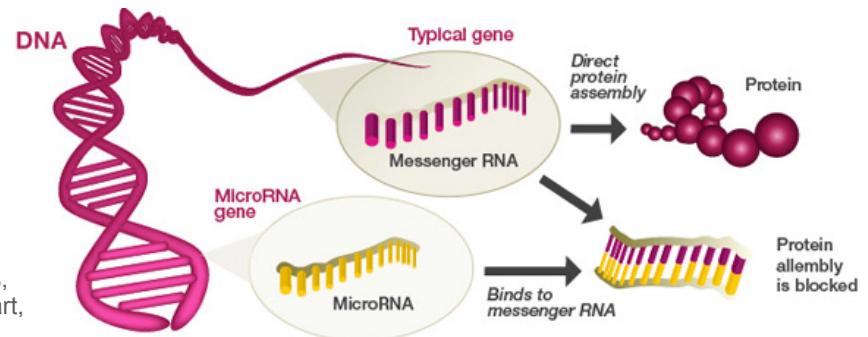
Multi-omics integration in cancer



RPPA → Gene expression
(Transcription Factors)

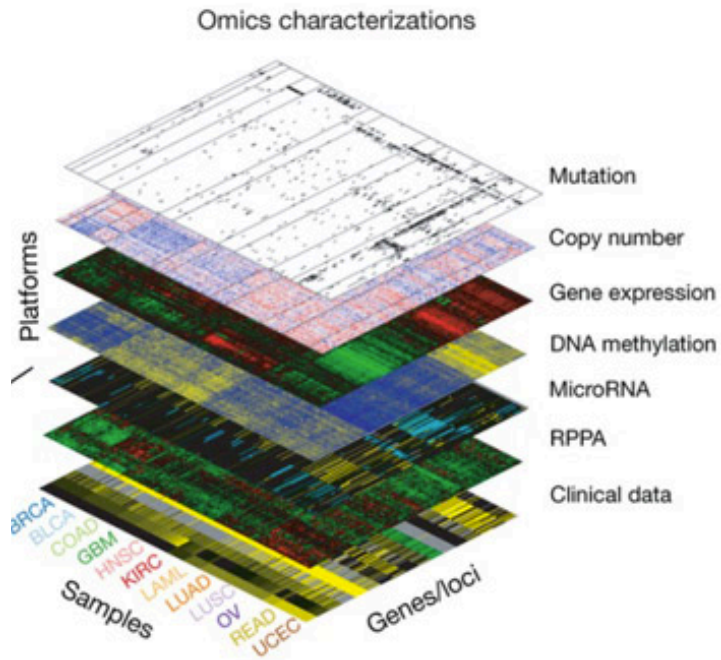


miRNA → RPPA
(microRNAs)



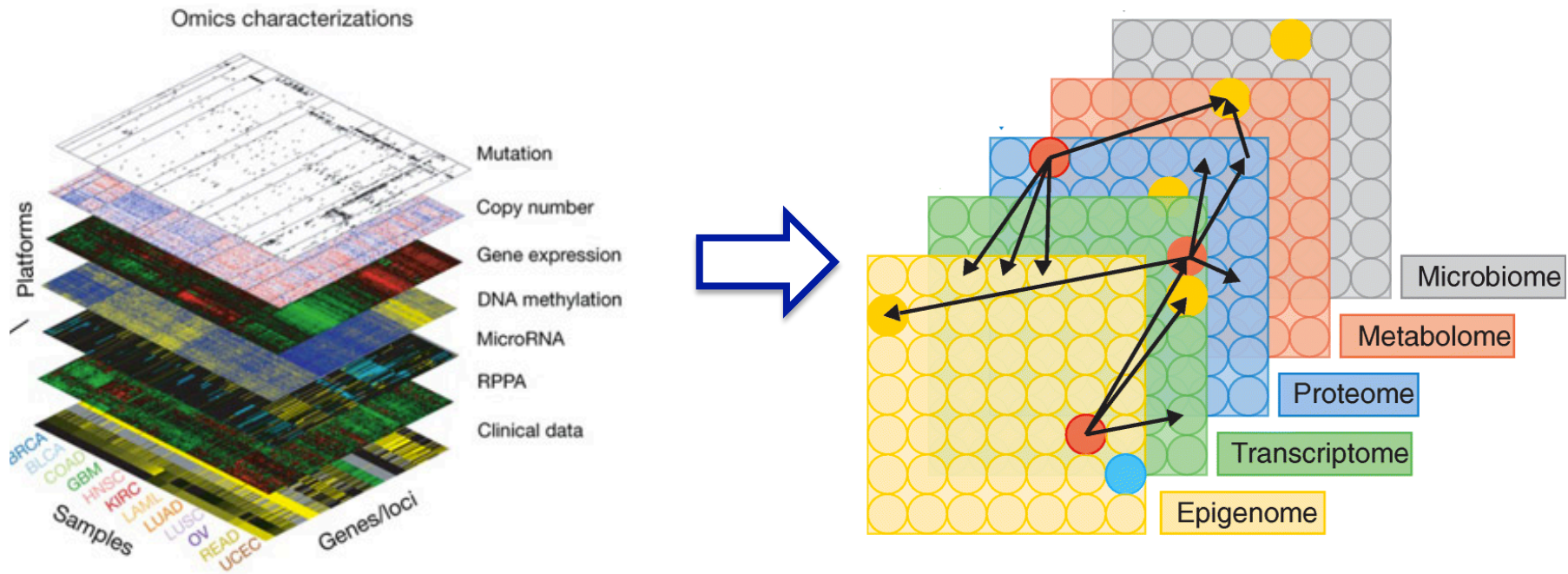
The Cancer Genome Atlas Research Network, Weinstein, J.N., Collisson, E.A., Mills, G.B., Shaw, K.M., Ozenberger, B.A., Ellrott, K., Shmulevich, I., Sander, C., and Stuart, J.M. (2013) The Cancer Genome Atlas Pan-Cancer analysis project. Nat Genet. doi: 10.1038/ng.2764

Multi-omics integration in cancer



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

Multi-omics integration in cancer



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

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



**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:

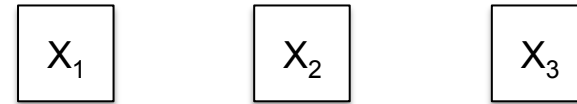
- Mutations
- Copy Number Alterations (CNA)
- Methylation
- Gene expression
- Proteomics
- Cancer types
- Drug response

Choosing which omics to integrate

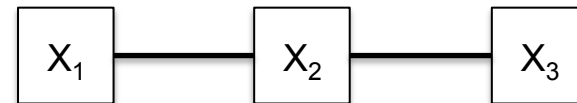
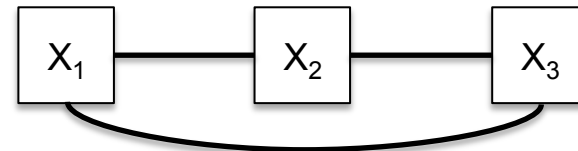
Aim: predicting drug response

Available input data:

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



Using correlation:

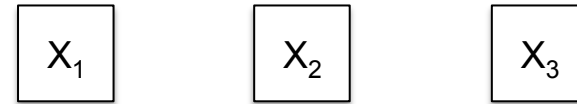


Choosing which omics to integrate

Aim: predicting drug response

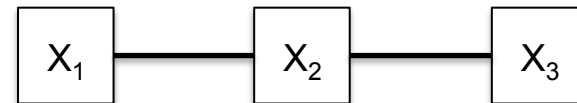
Available input data:

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

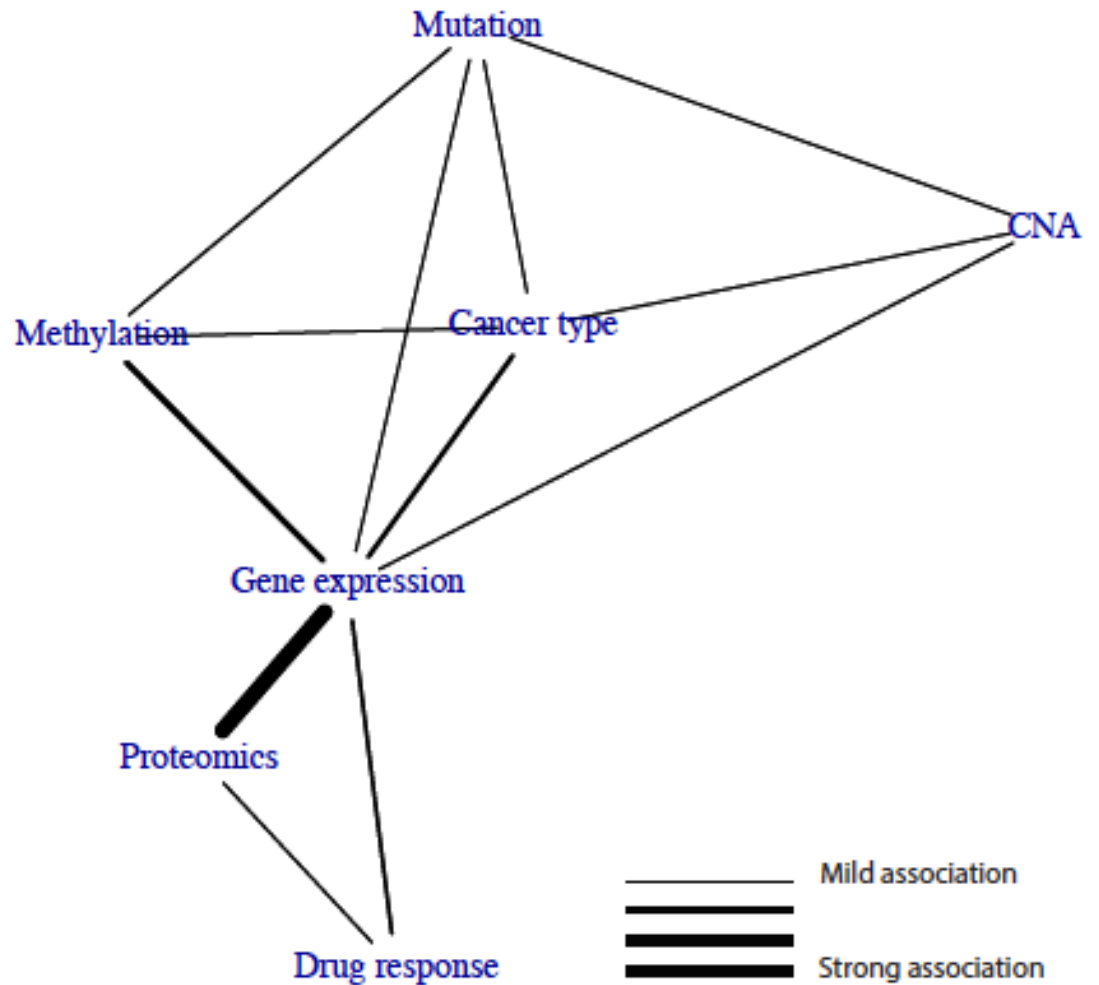


Using partial correlation (iTOP):

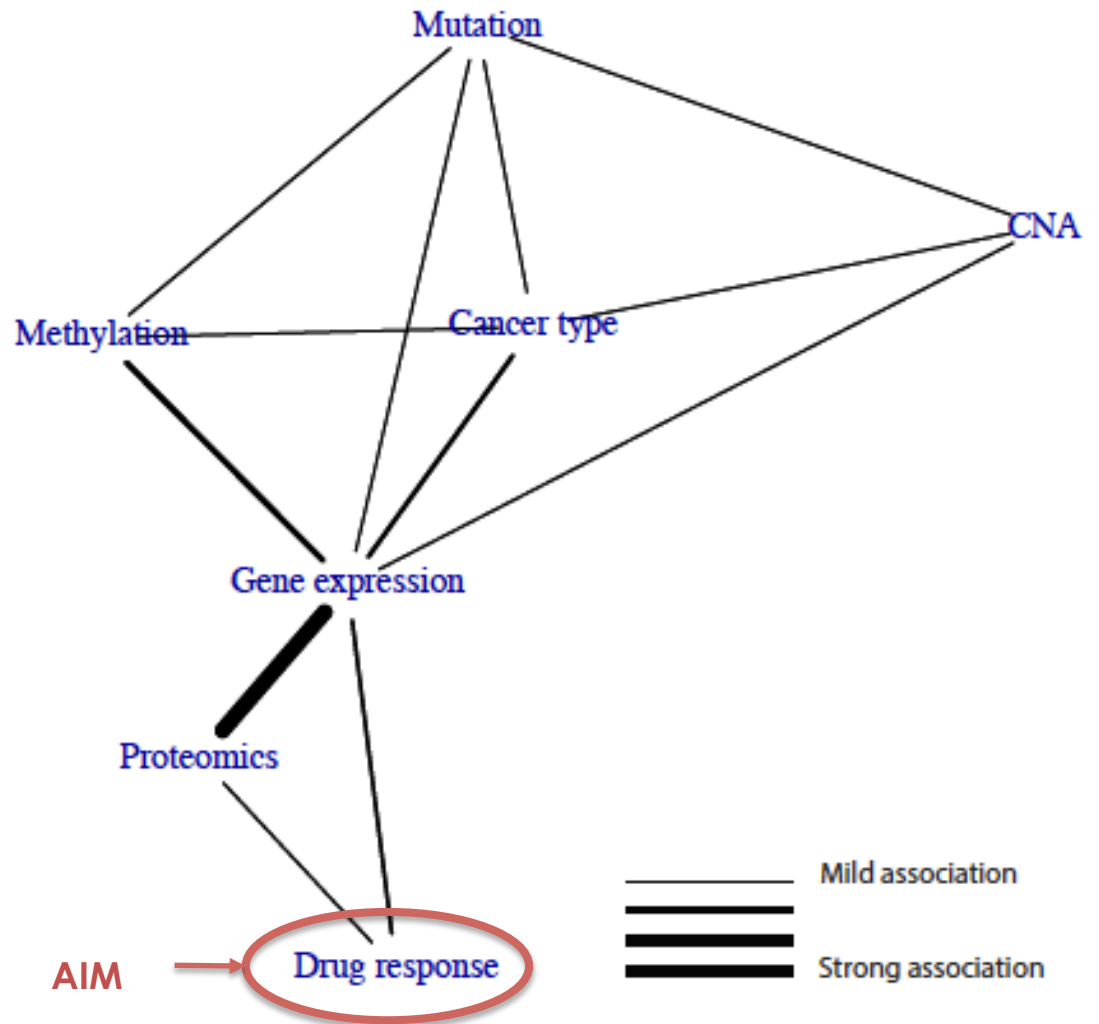
e.g. $\text{cor}(X_1, X_3 | X_2) \approx 0$



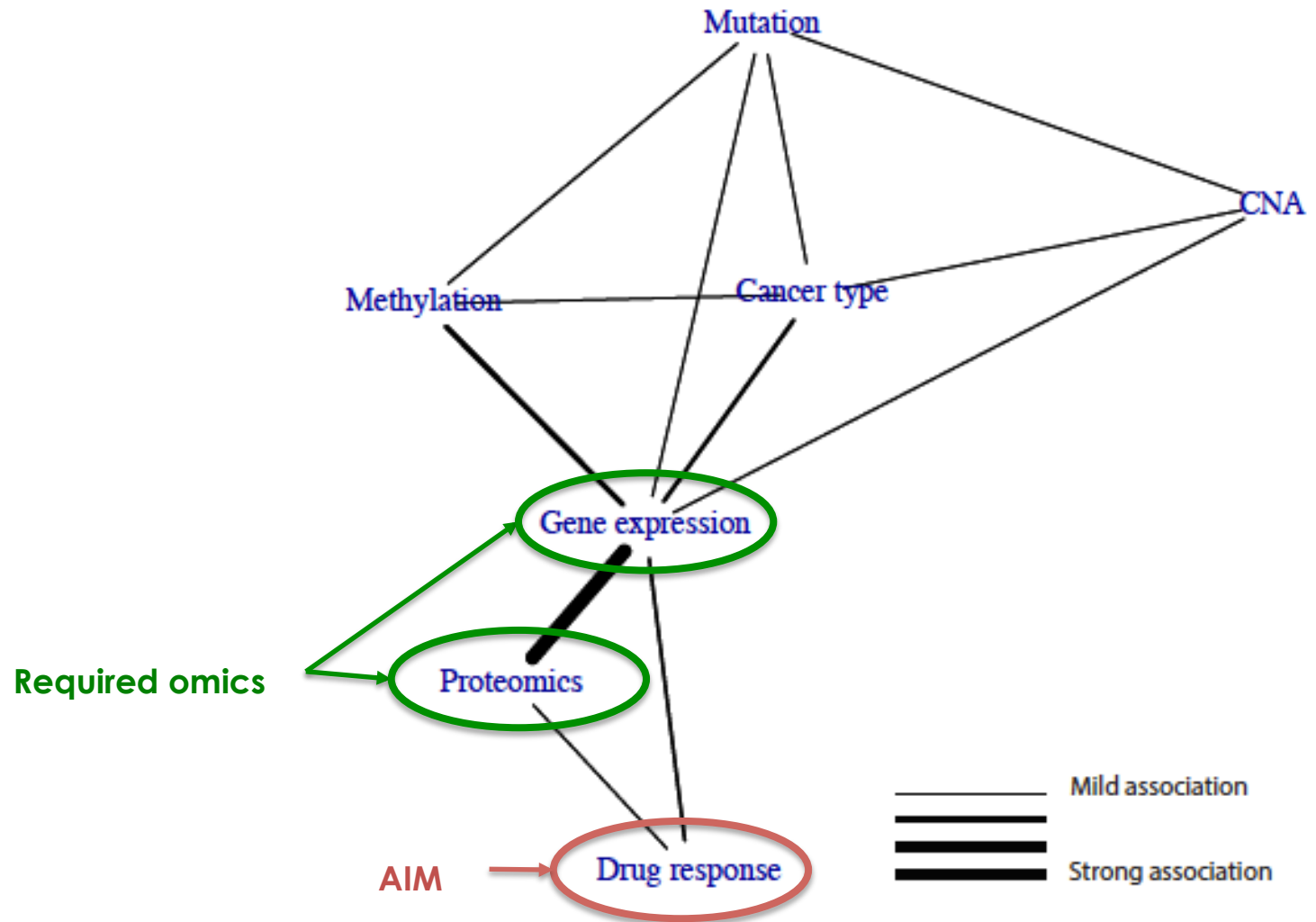
Choosing which omics to integrate



Choosing which omics to integrate



Choosing which omics to integrate

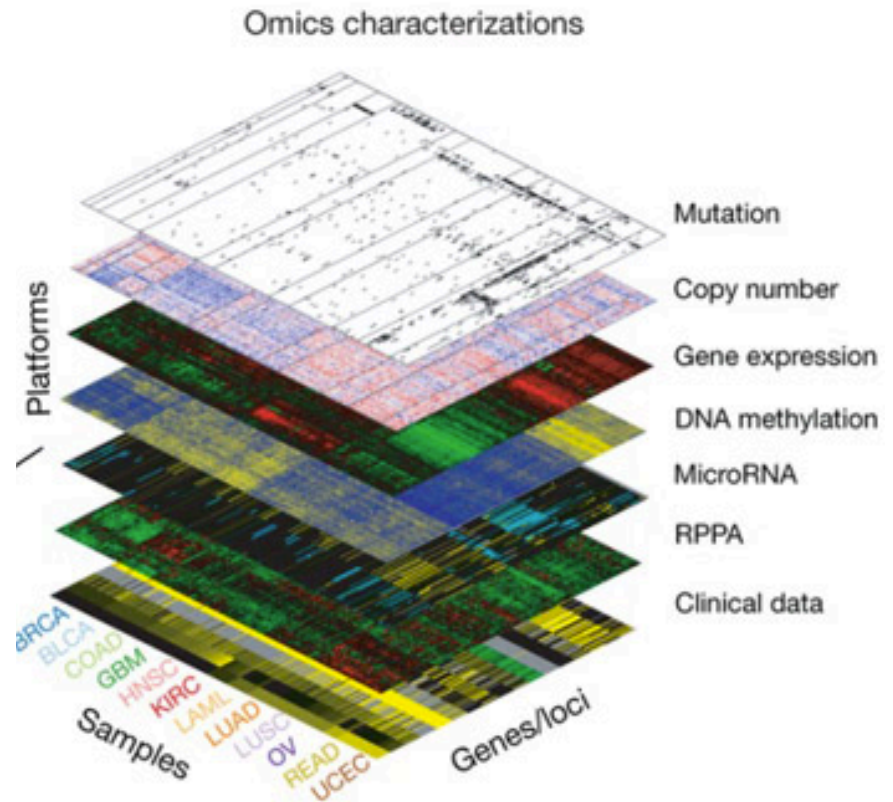




How omics should be combined?

Challenge: Multi-omics data are heterogeneous

- continuous (e.g. gene expression) vs. discrete (e.g. CNV).
- They have different dimension
- They have different ranges of variability
- Involving heterogeneous entities (e.g. genes, proteins, CpGs). -> matching entities drastically reduce the information

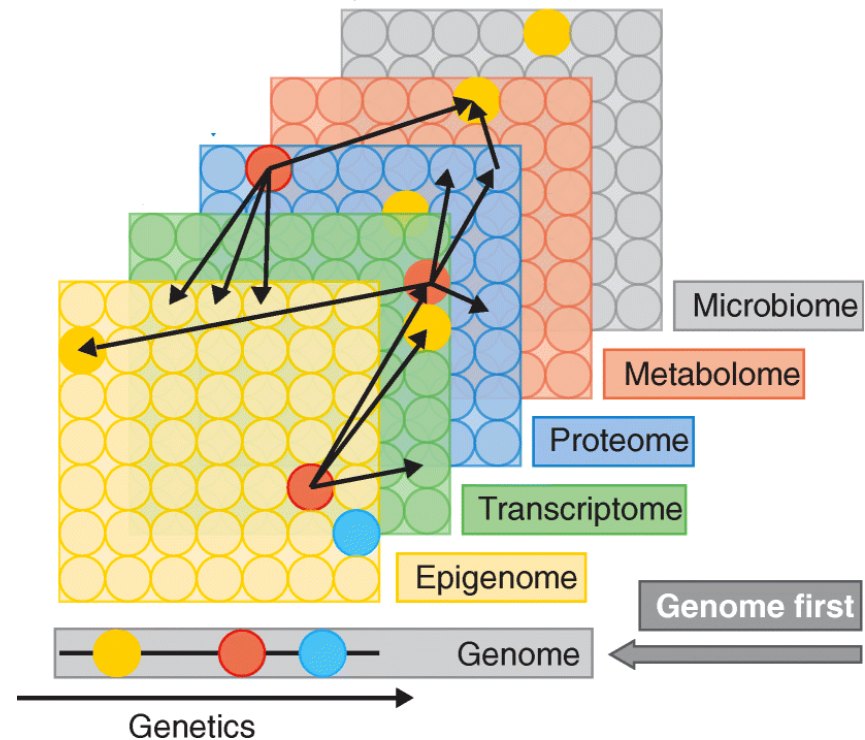


Integrating multi-omics data

Approach “Genome First”

Priority given to genome

Other omics are only used for interpretation

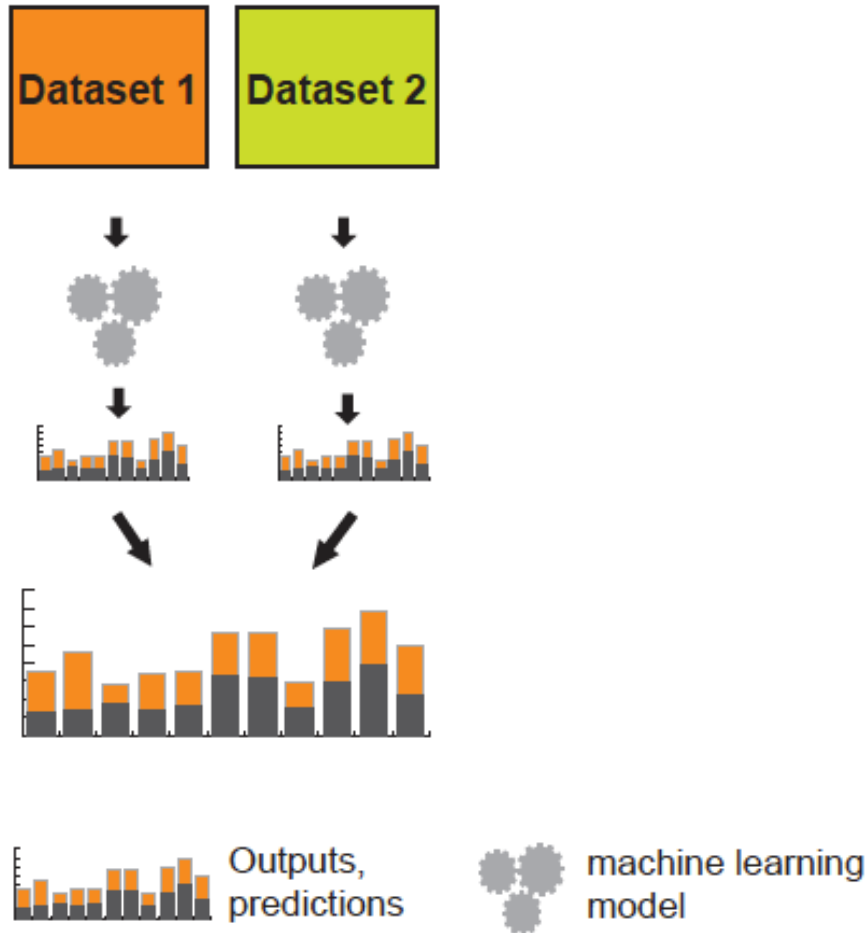


Hasin, Yehudit, Marcus Seldin, and Aldons LUIS. "Multi-omics approaches to disease." *Genome biology* 18.1 (2017): 83.

Integrating multi-omics data

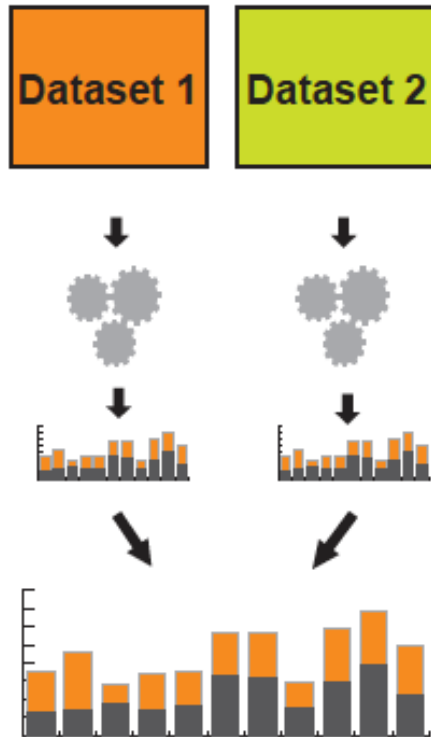
Late integration

output averaging, ensembles

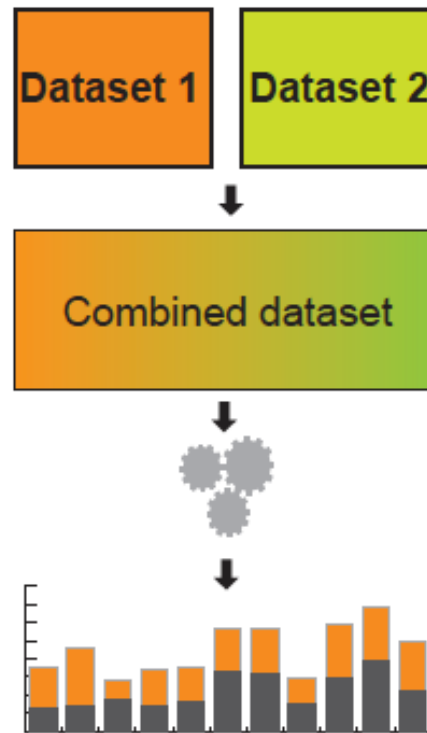


Integrating multi-omics data


Late integration
output averaging, ensembles



Early integration
projection, concatenation

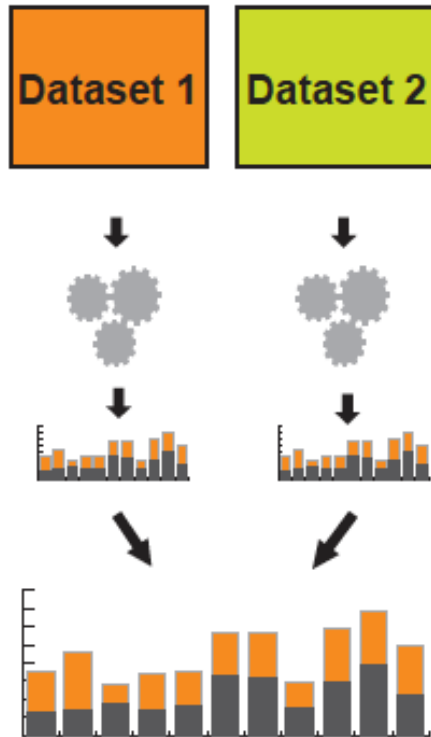


 Outputs,
predictions

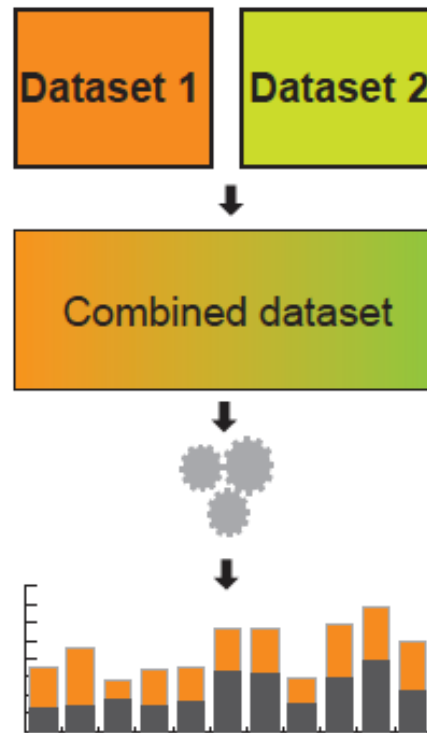
 machine learning
model

Integrating multi-omics data

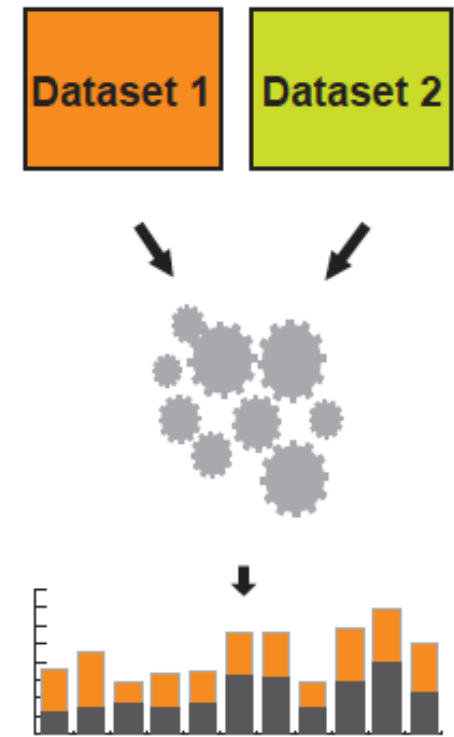
Late integration
output averaging, ensembles




Early integration
projection, concatenation



Intermediate integration
multi-view, multi-modal



 Outputs,
predictions

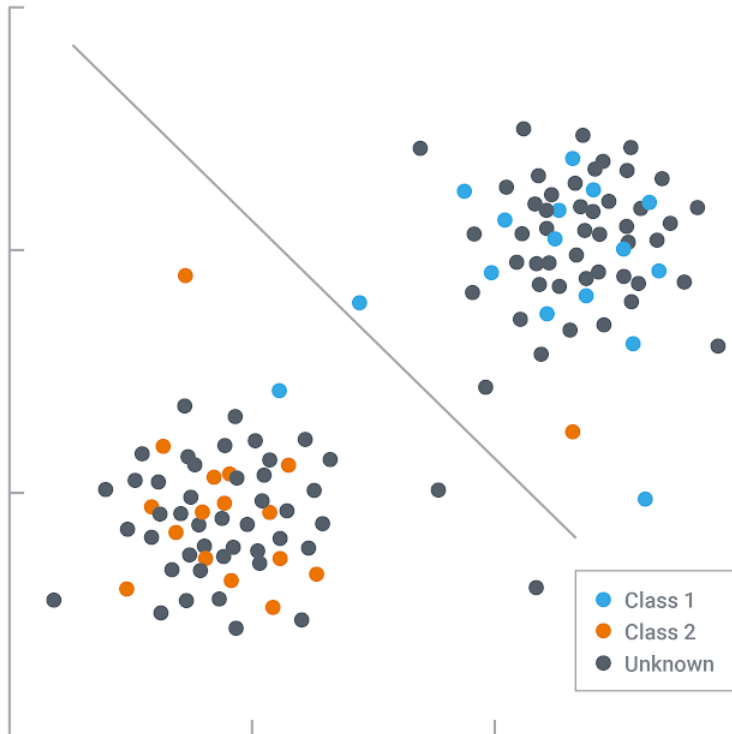
 machine learning
model



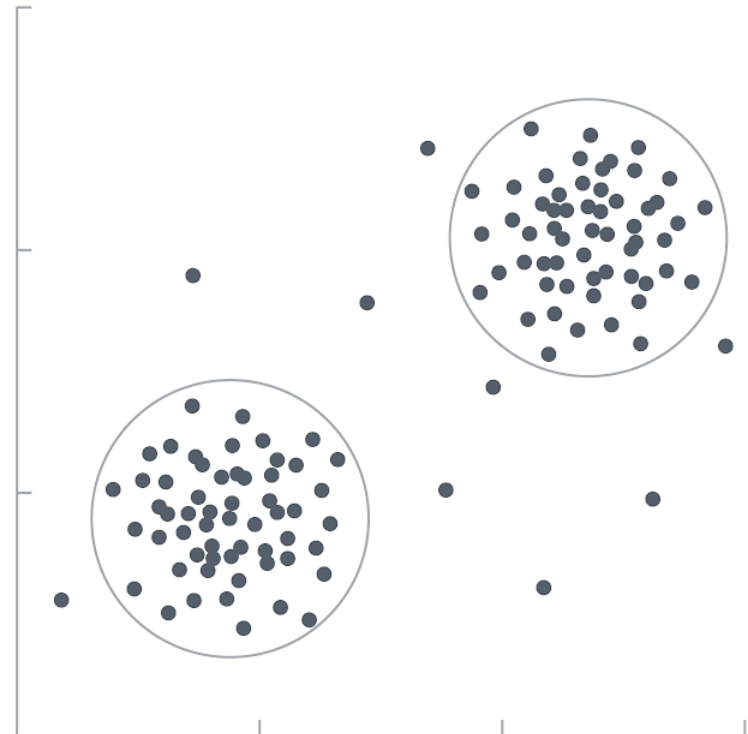
**Main categories of
existing multi-omics
integrative approaches**

Main categories of integrative approaches

Supervised methods

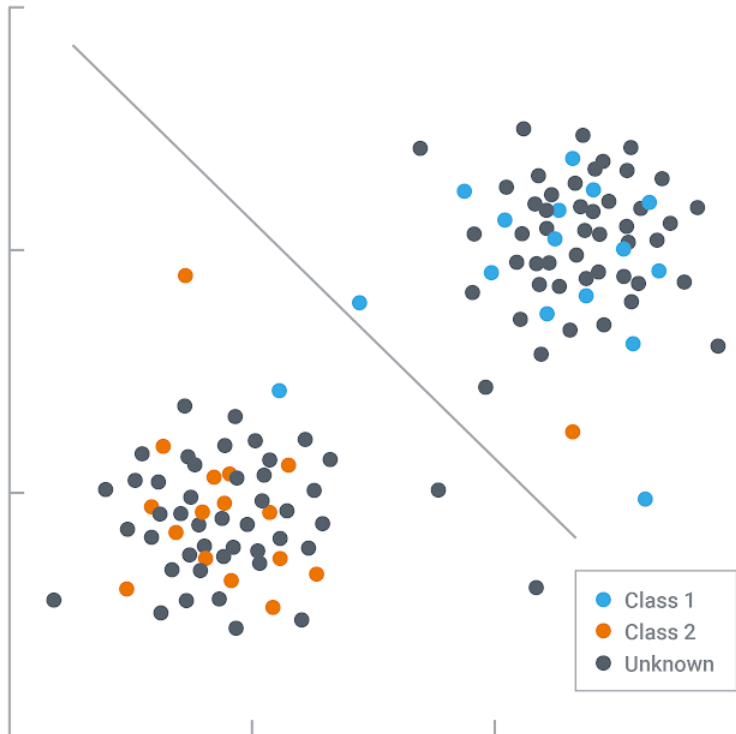


Unsupervised methods



Main categories of integrative approaches

Supervised methods

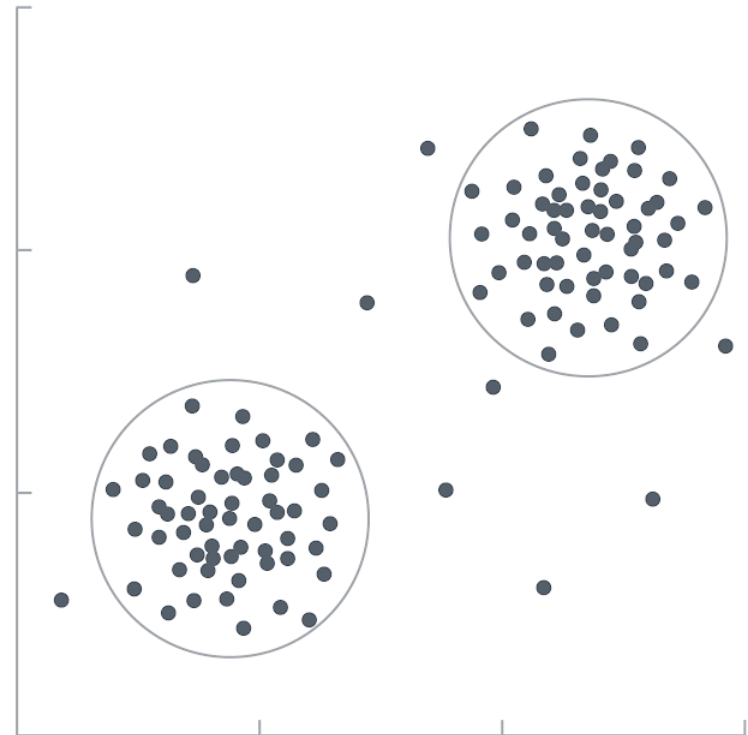


- They require 2 datasets in input: training and test datasets
- Labels must be available for the training dataset
- This information is used to infer labels on the test dataset

Main categories of integrative approaches

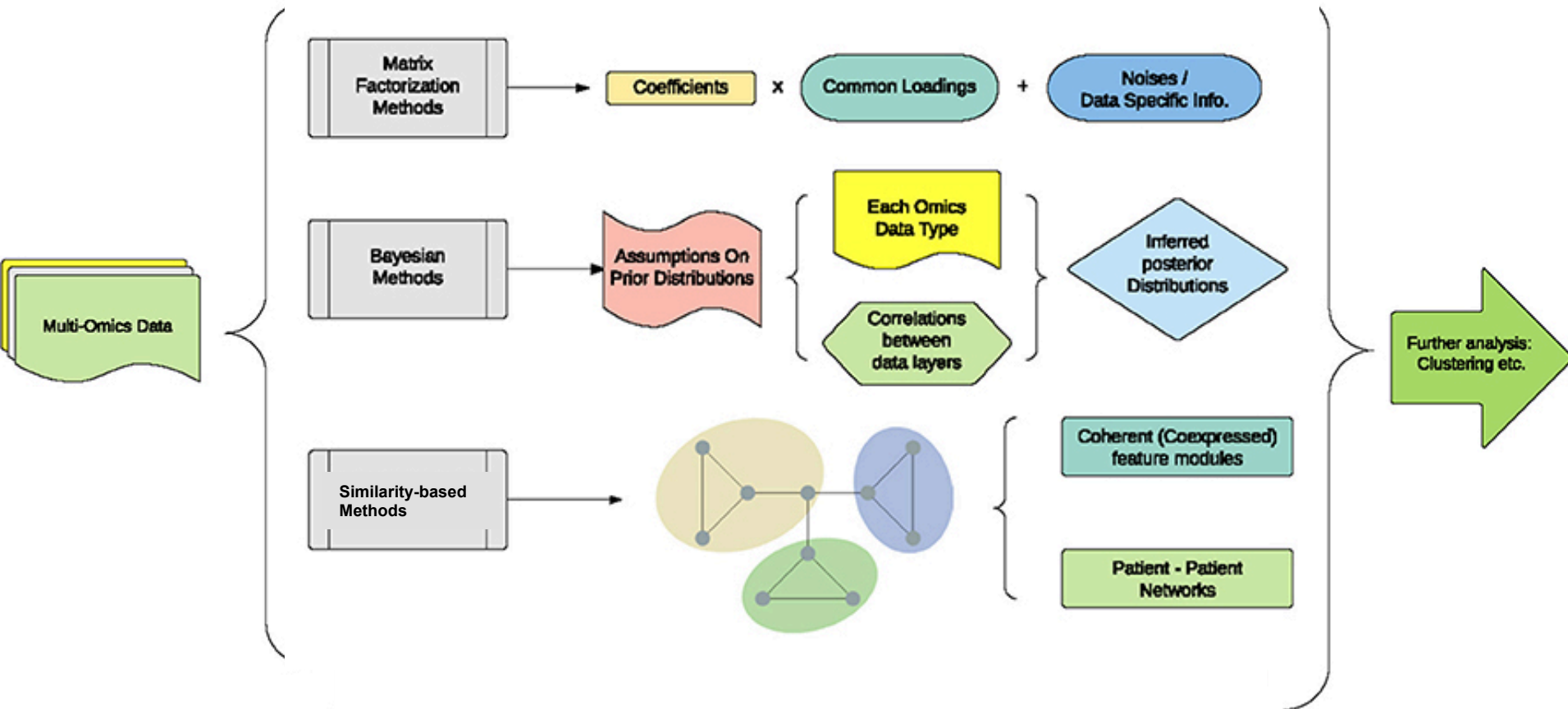
- The methodology is directly applied to one dataset
- They infer information from the structure of the data without any label information

Unsupervised methods



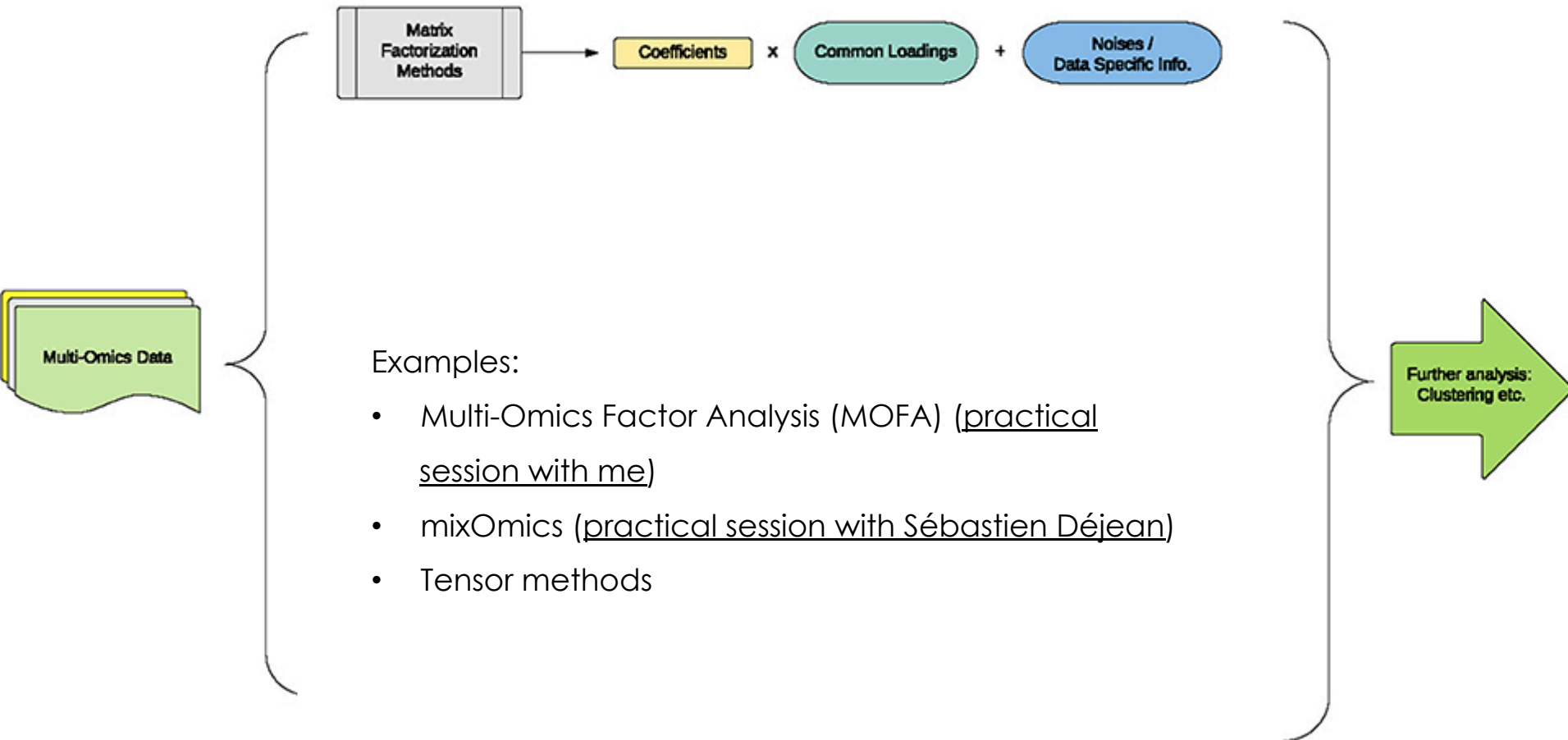
Unsupervised integrative approaches

Unsupervised data integration



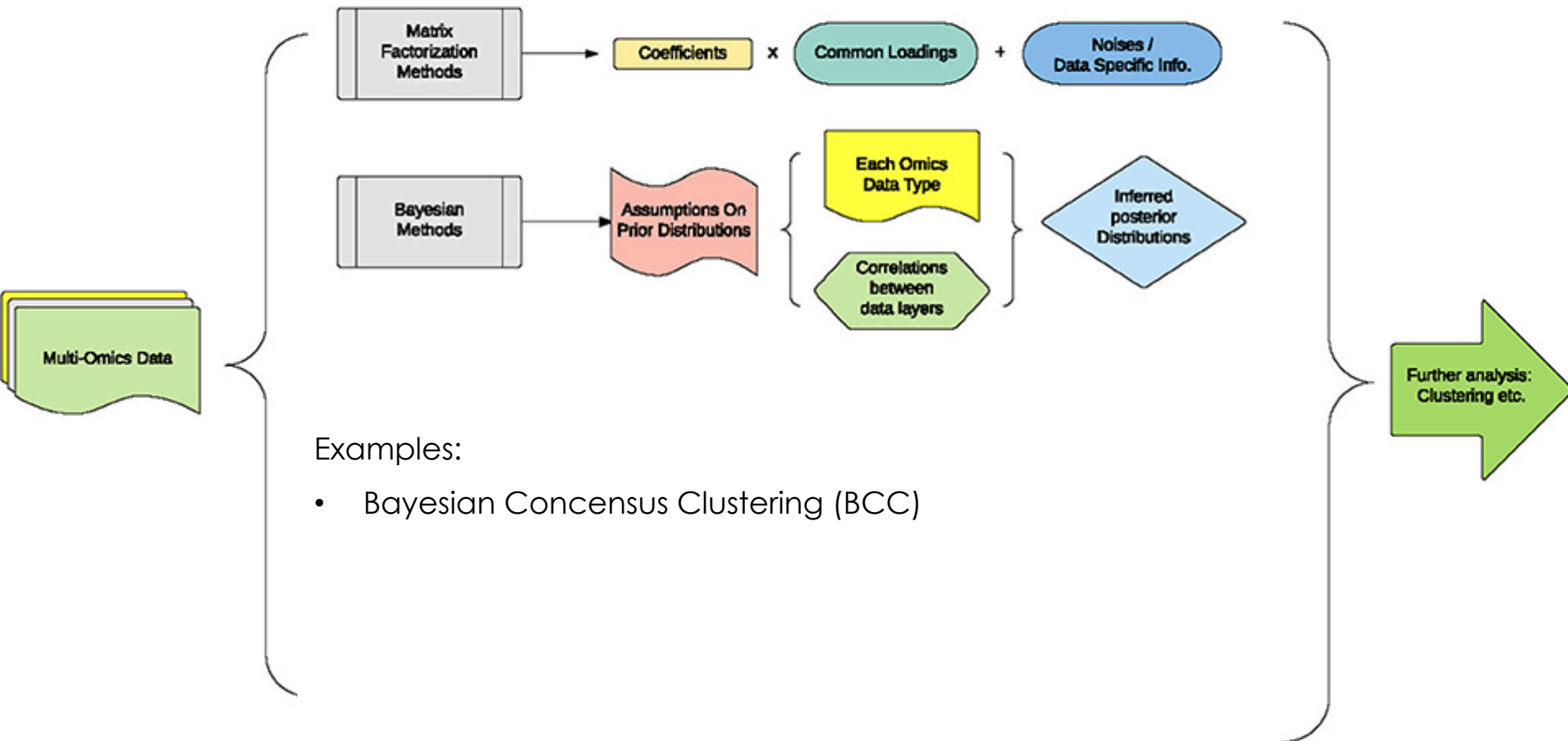
Unsupervised integrative approaches

Unsupervised data integration



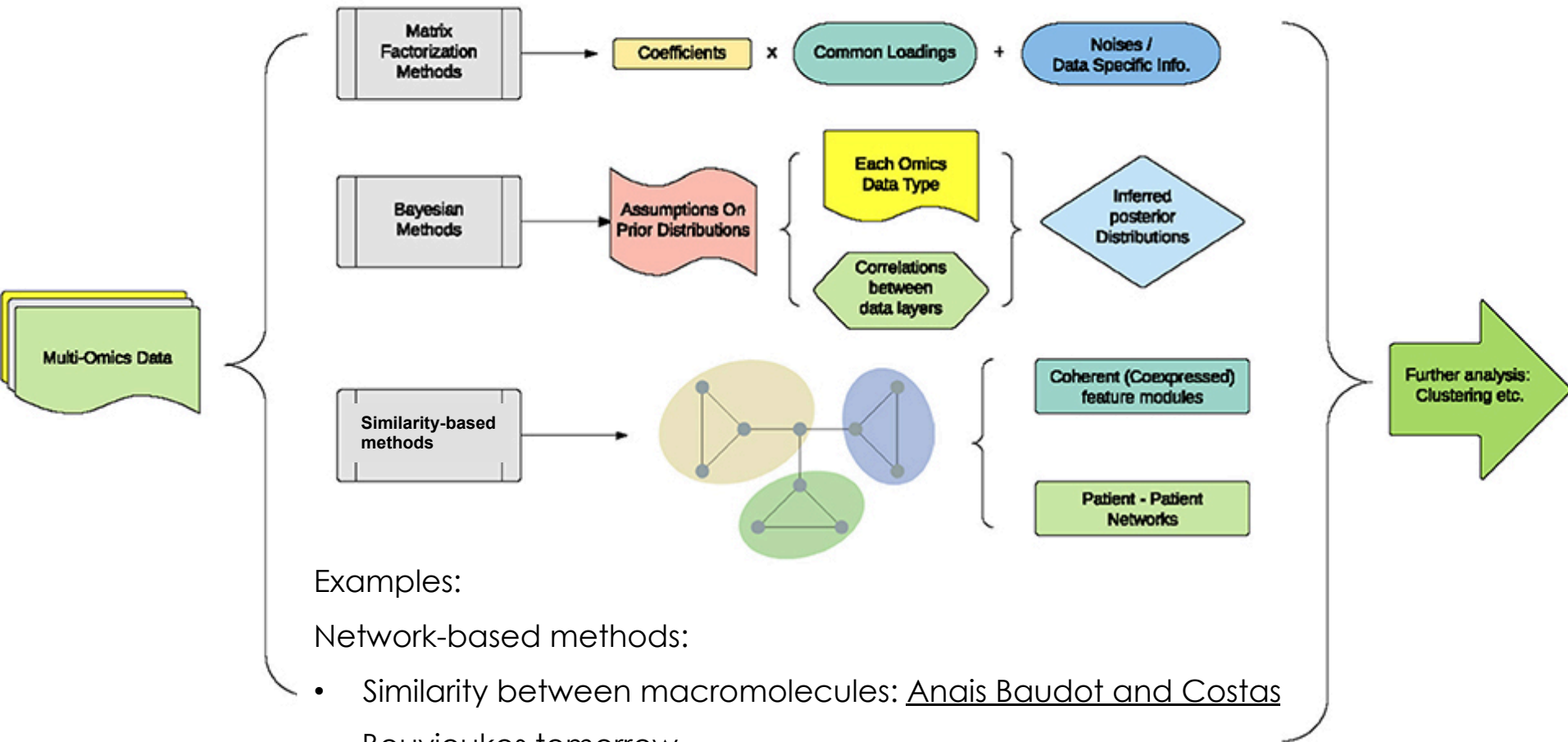
Unsupervised integrative approaches

Unsupervised data integration



Unsupervised integrative approaches

Unsupervised data integration



Examples:

Network-based methods:

- Similarity between macromolecules: Anais Baudot and Costas Bouyioukos tomorrow
- Similarity between samples: Similarity Network Fusion SNF

Kernel methods: practical session Jérôme Mariette

Unsupervised integrative approaches

Of Note:

This last category of methods give the possibility to integrate heterogeneous sources of information, such as: data matrices, phylogenetic trees and networks

	Sample_a	Sample_b	Sample_c
OTU_1	3	3	3
OTU_2	5	2	3
OTU_3	2	0	2
OTU_4	4	3	0
OTU_5	0	2	2
OTU_6	0	2	2



Examples:

Network-based methods:

- Similarity between macromolecules: Anais Baudot tomorrow
- Similarity between samples: Similarity Network Fusion SNF

Kernel methods: practical session Jérôme Mariette

ysis:
etc.

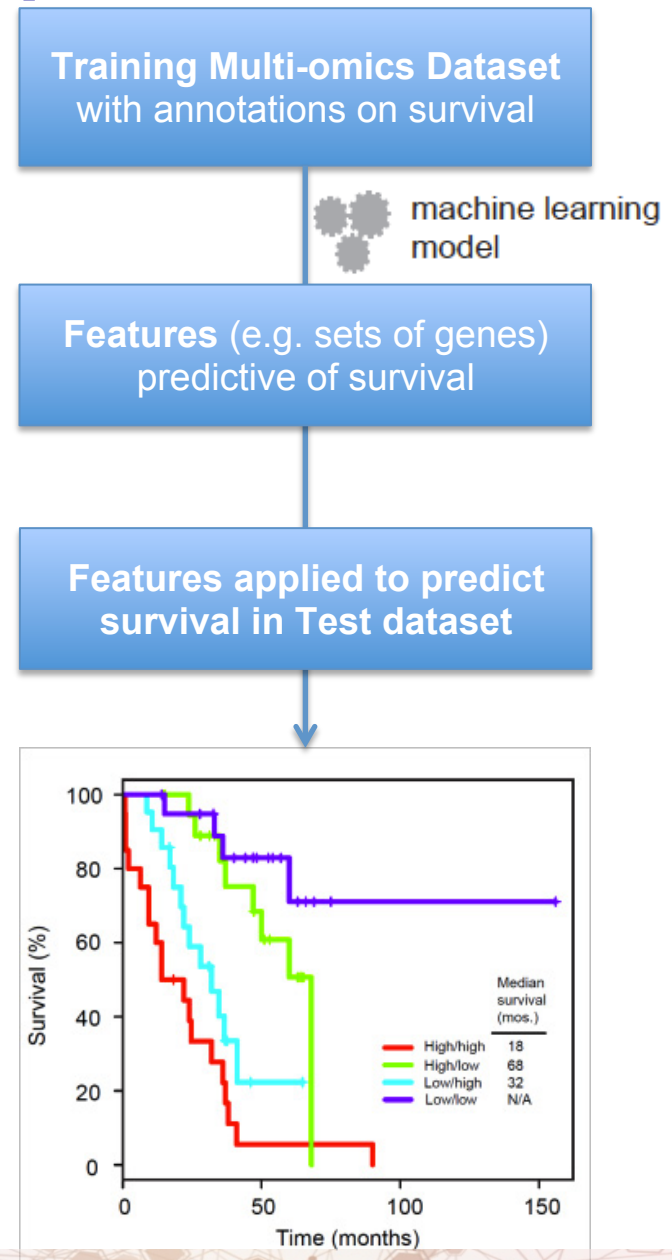


Cancer insights from data integration methods

Patients survival prediction

Given a set of cancer we want to patients predict their survival.

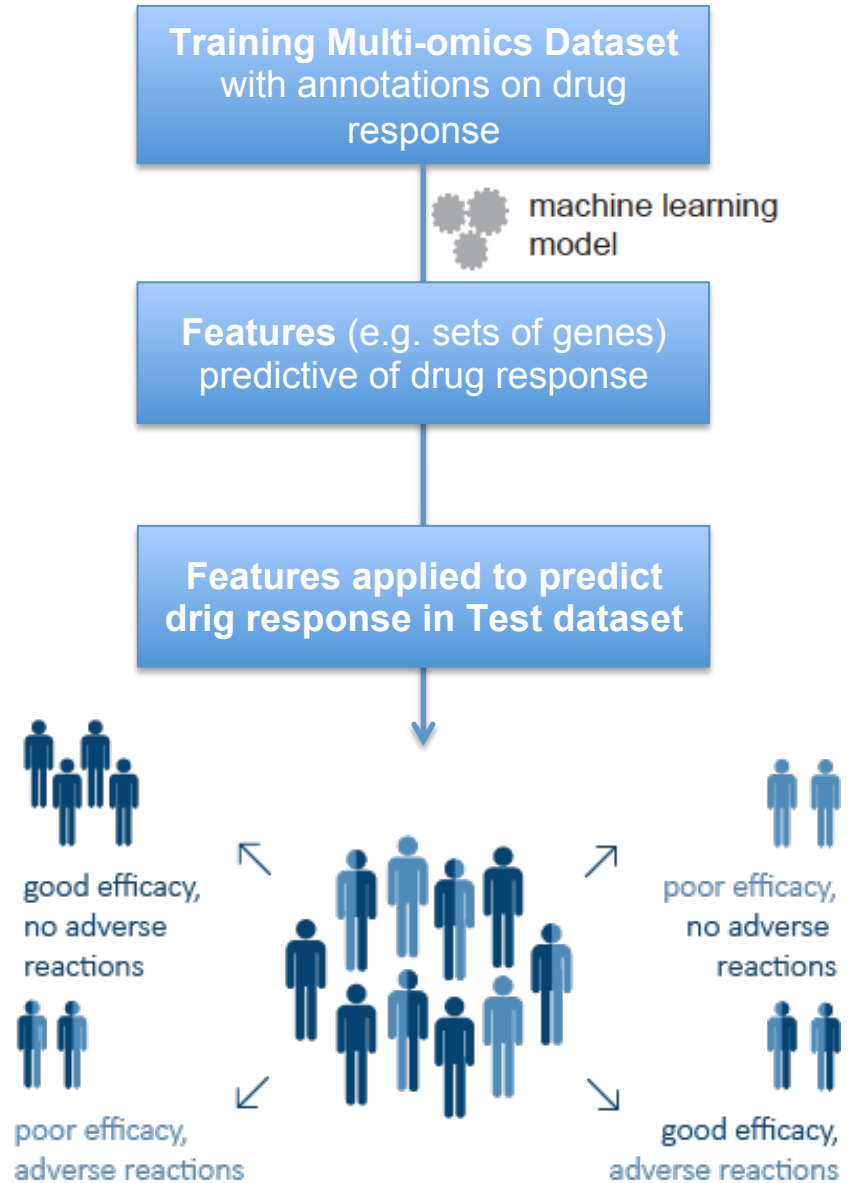
This problem is generally approached with supervised approaches.



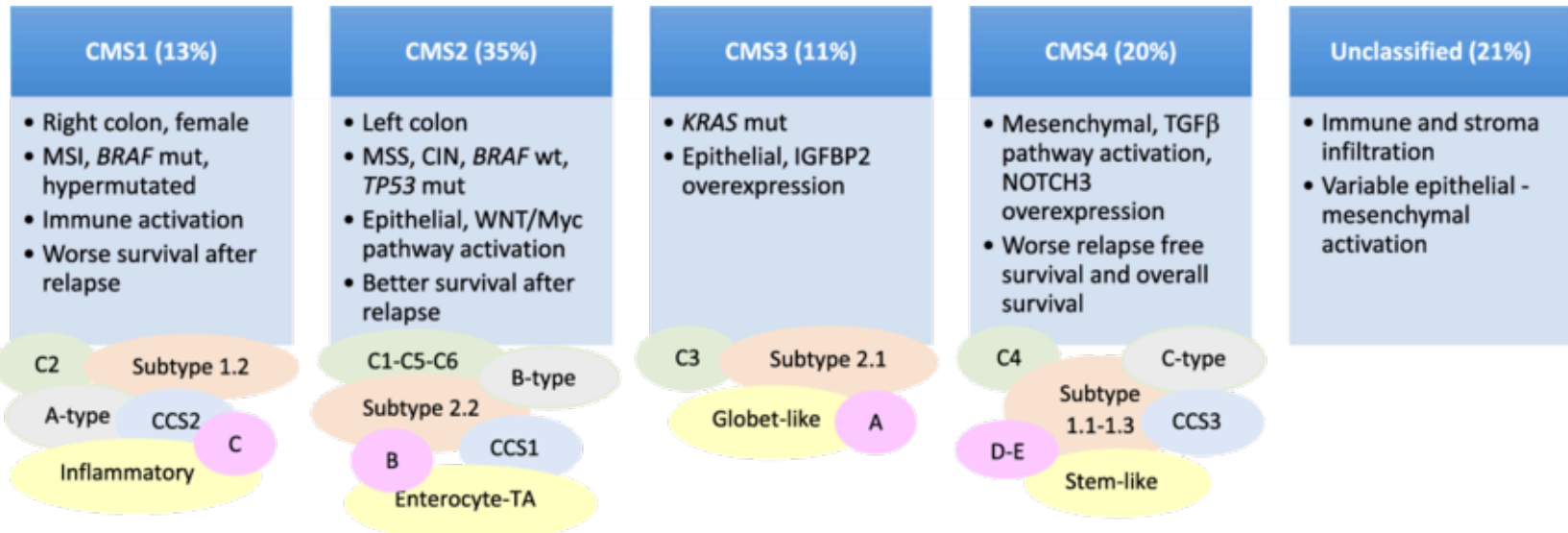
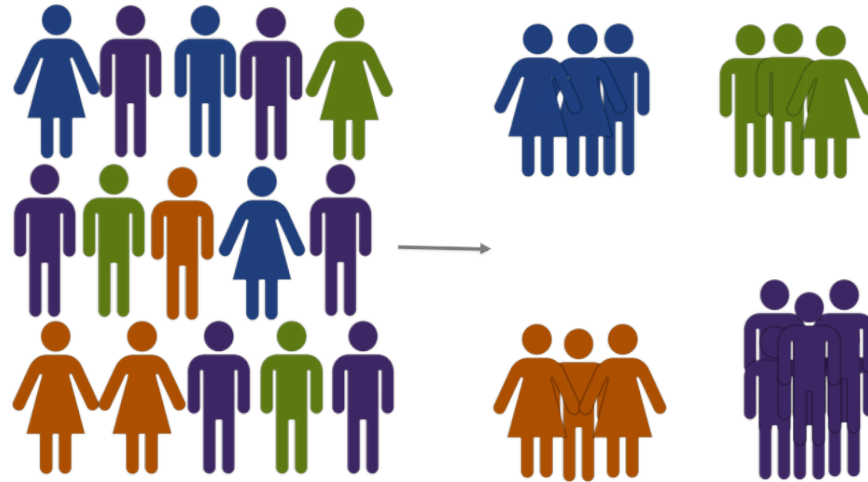
Drug response prediction

Given a set of cancer we want to patients predict which will respond to a given therapy.

This problem is generally approached with supervised approaches.



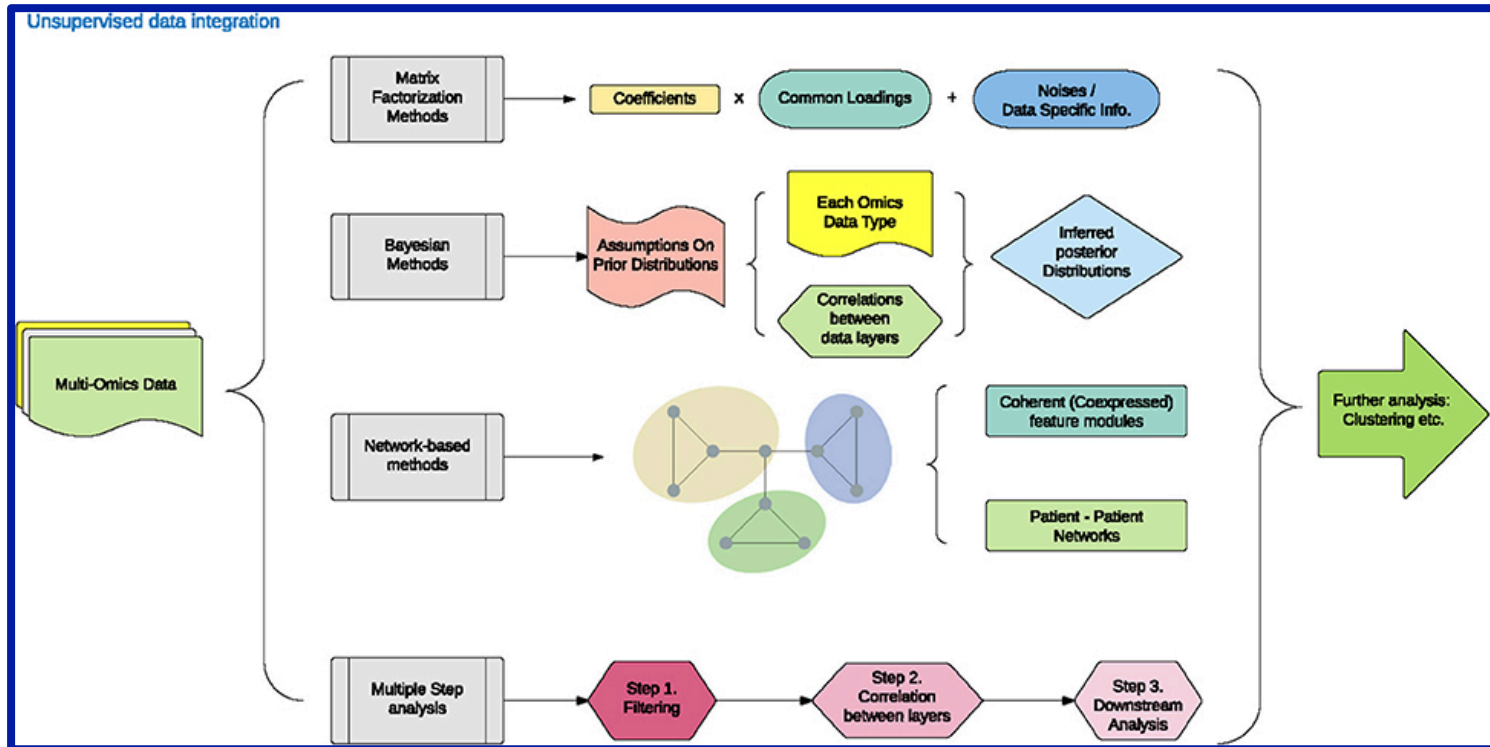
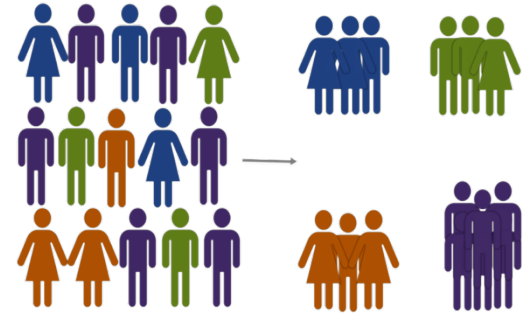
Cancer subtyping



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 responding



Drug resistant



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?

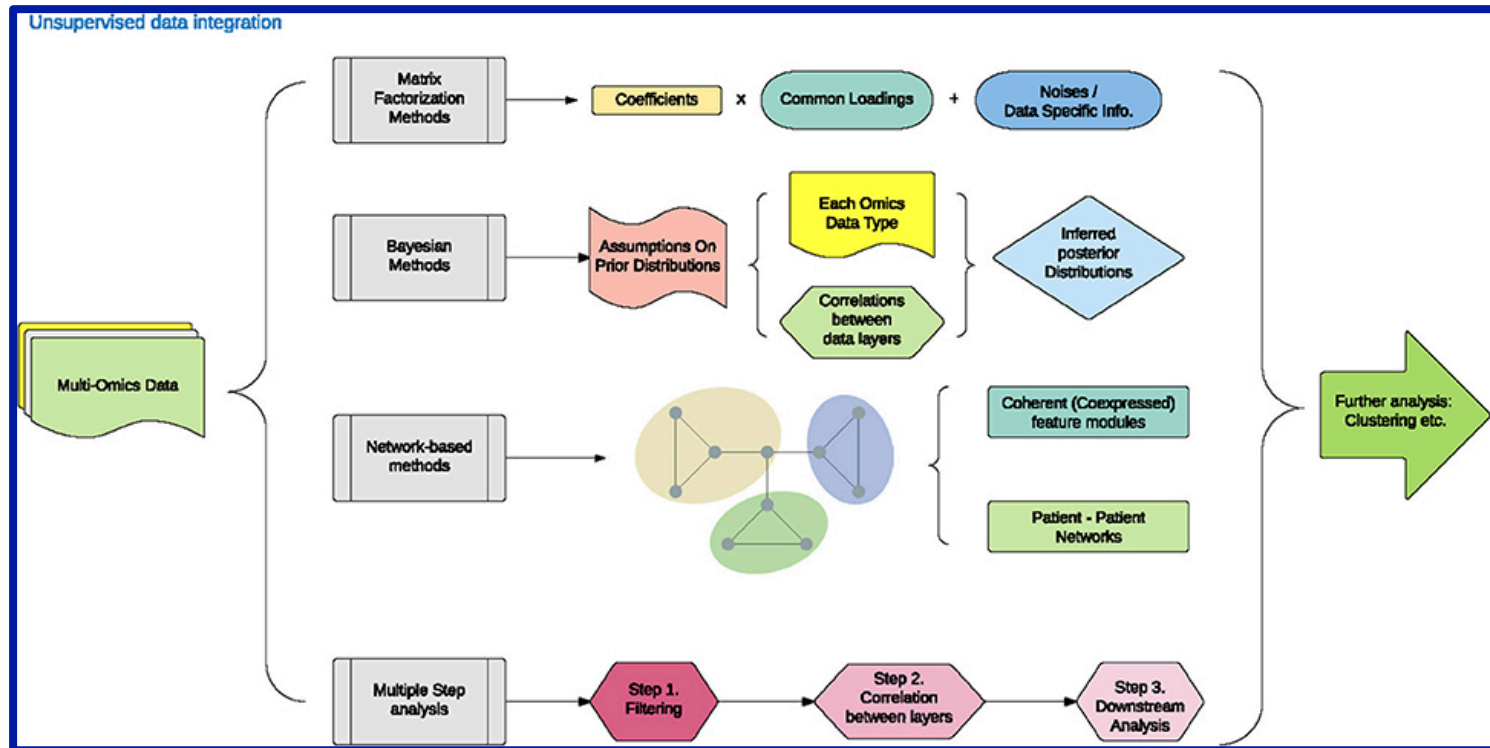
Gene modules identification

This problem is generally approached with unsupervised approaches.

Drug responding



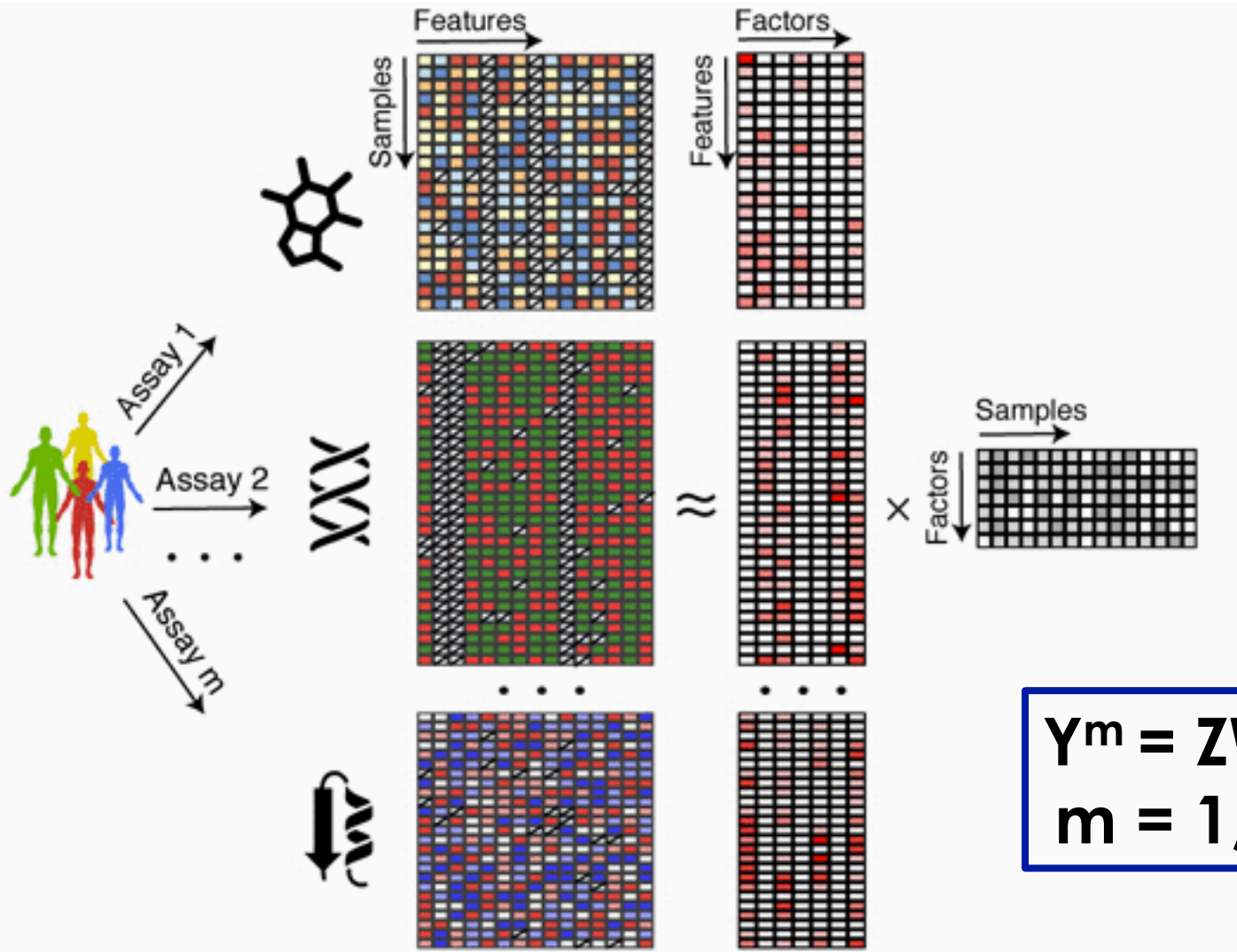
Drug resistant





Multi-Omics Factor Analysis (MOFA)

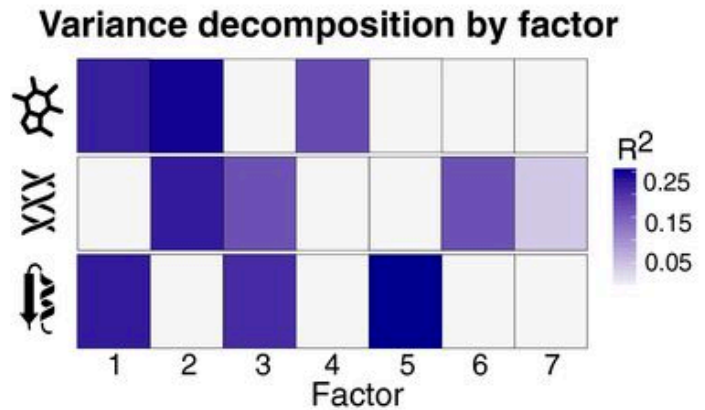
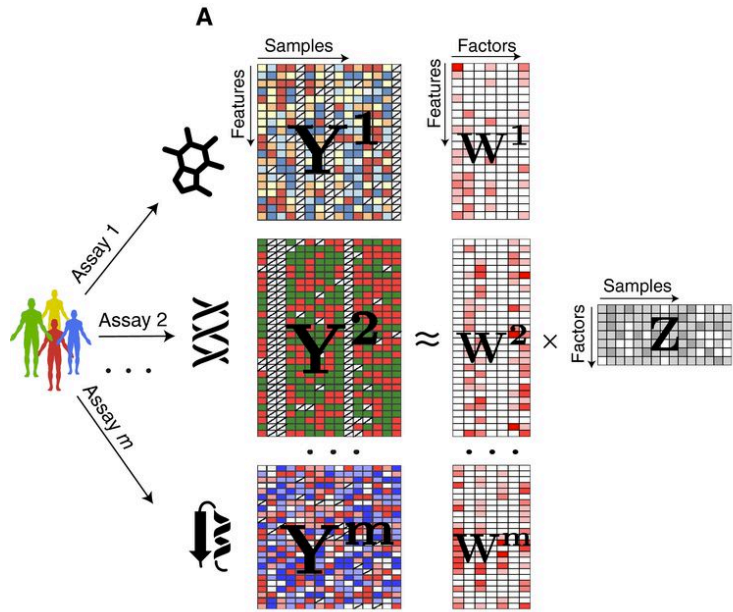
MOFA model



$$Y^m = ZW^m + e^m$$

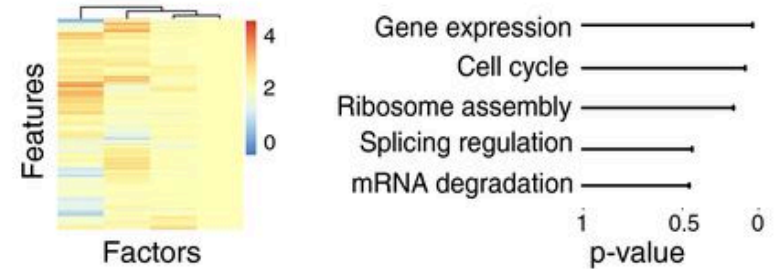
$$m = 1, \dots, M$$

MOFA advantage: interpretability of factors

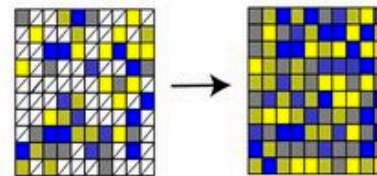


Annotation of factors

Inspection of loadings Feature set enrichment analysis



Imputation of missing values



Inspection of factors

