

Intégration Multi-omiques

Antoine Bodein

Charles Joly-Beauparlant

Arnaud Droit

Paris Diderot 20-21 octobre 2020

© Département de biochimie, de microbiologie et de bio-informatique, 2020



UNIVERSITÉ
Laval

Objectifs

- ❑ Comprendre les **enjeux** liés à l'intégration multi-omique.
- ❑ Être capable d'expliquer les différents **concepts** d'intégration multi-omique.
- ❑ Savoir appliquer la bonne **méthode** d'intégration en fonction des **données** disponibles.

Mise en contexte

- omique ? multi-omique ?
- Pourquoi ? Difficultés
- Ressources

Méthodes d'intégration

- Intégration Conceptuelle
- Intégration Statistique
- Intégration à base de modèle

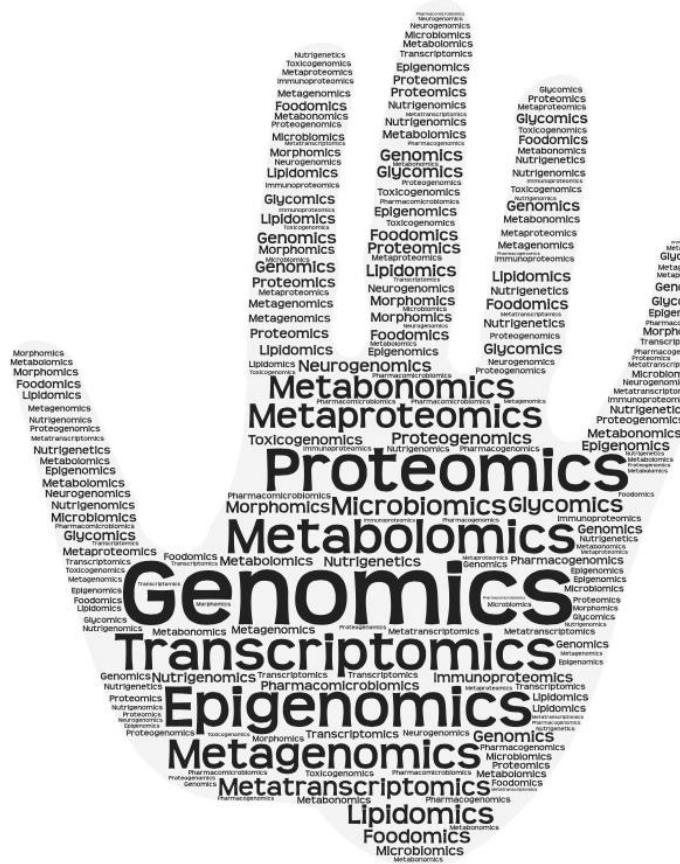
Méthodes Multivariées

- Analyse en Composante Principale
- Autres méthodes multivariées

Méthodes à base de réseaux biologiques

- Les réseaux en biologie
- Exploration des réseaux

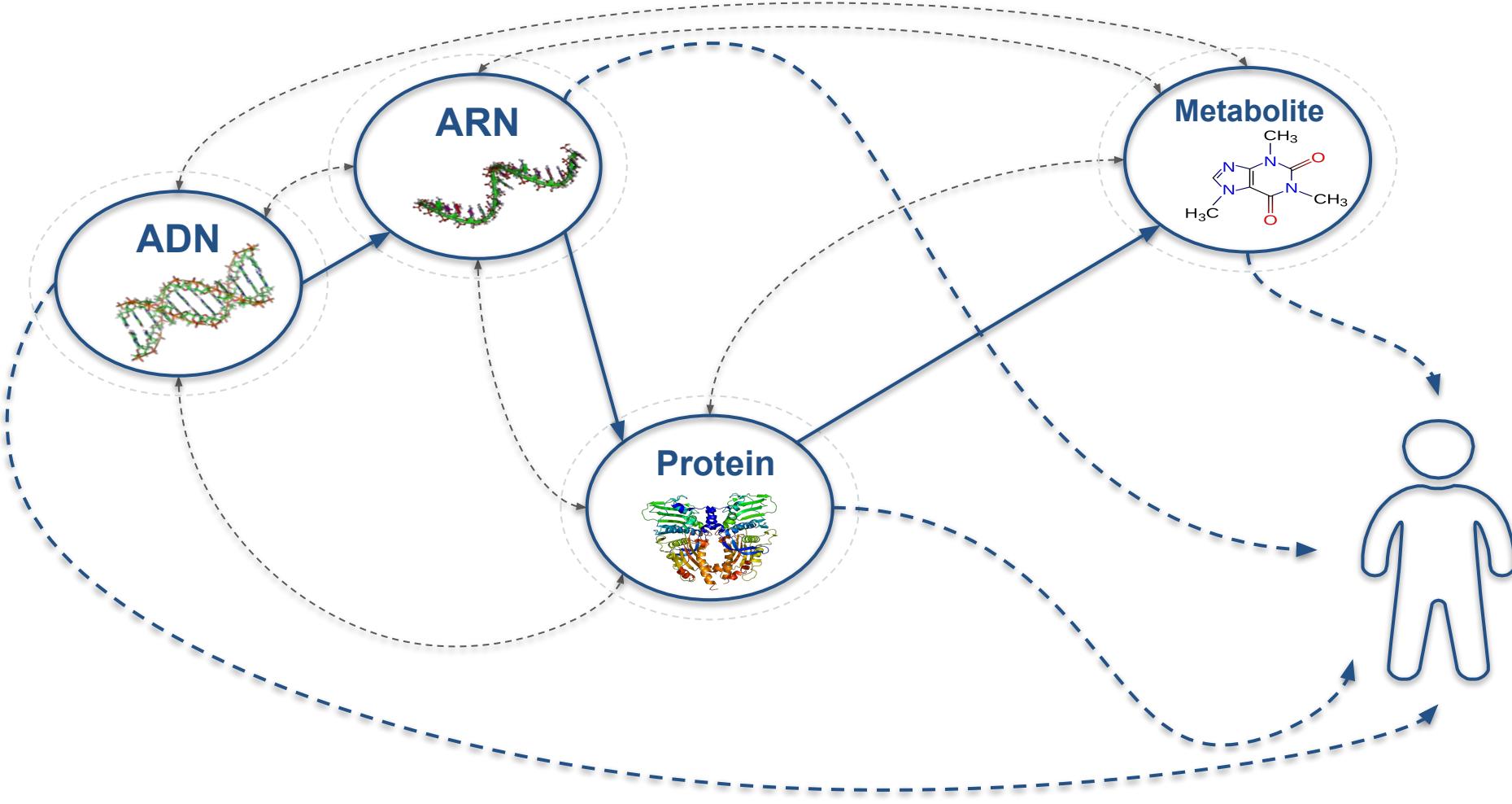
Qu'est-ce qu'un *omique* ?



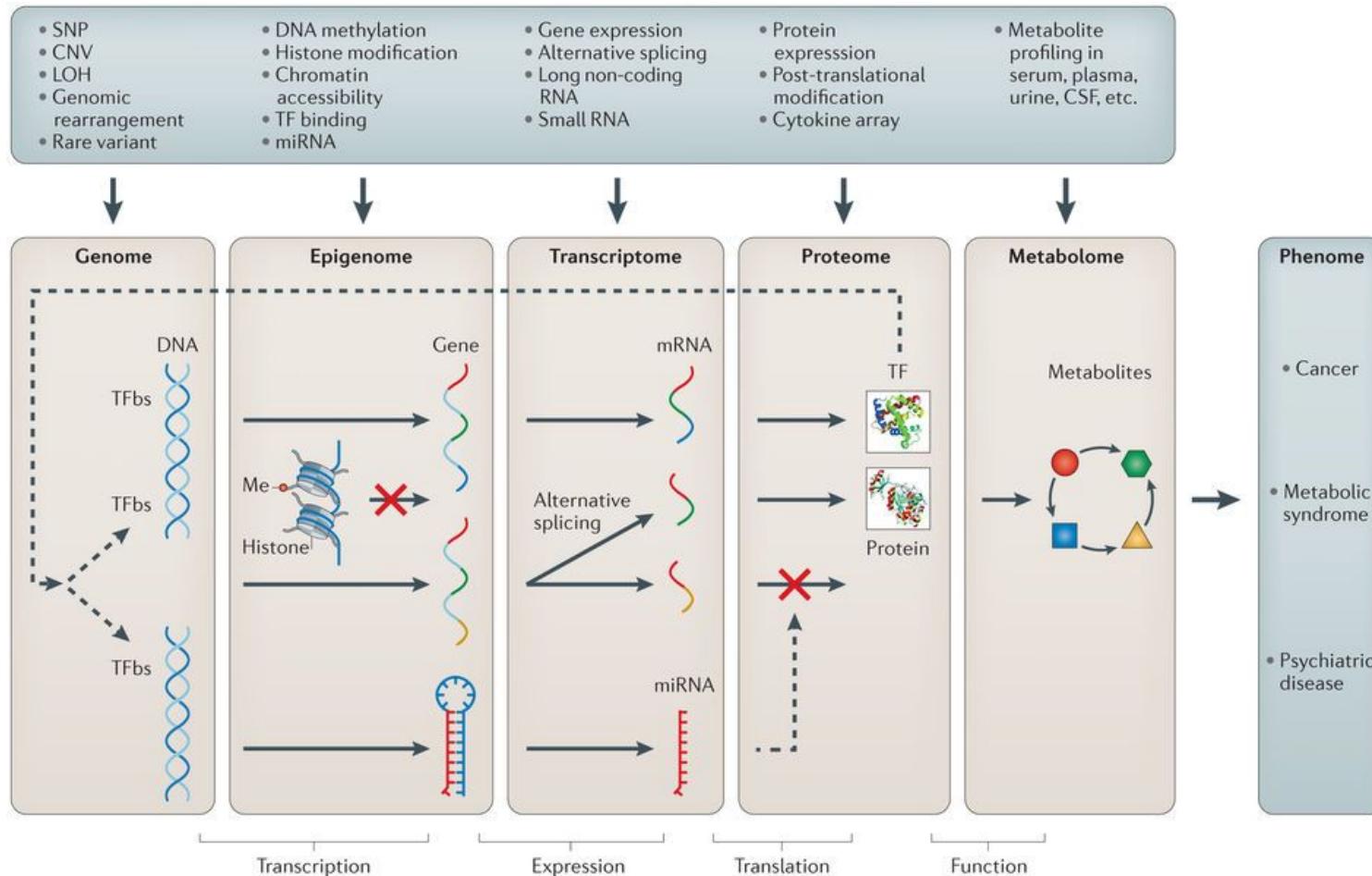
- ❑ “Genomics” inventé dans les 70s
(Thomas Roderick)
 - ❑ D’autres termes:
Protéomique, métabolomique, épigénomique, ...
 - ❑ “Méta”-omique
Metagénomique, métaprotéomique, ...
 - ❑ Néologisme
Morphomics, foodomics, ...

“Technologies omiques à haut débit”

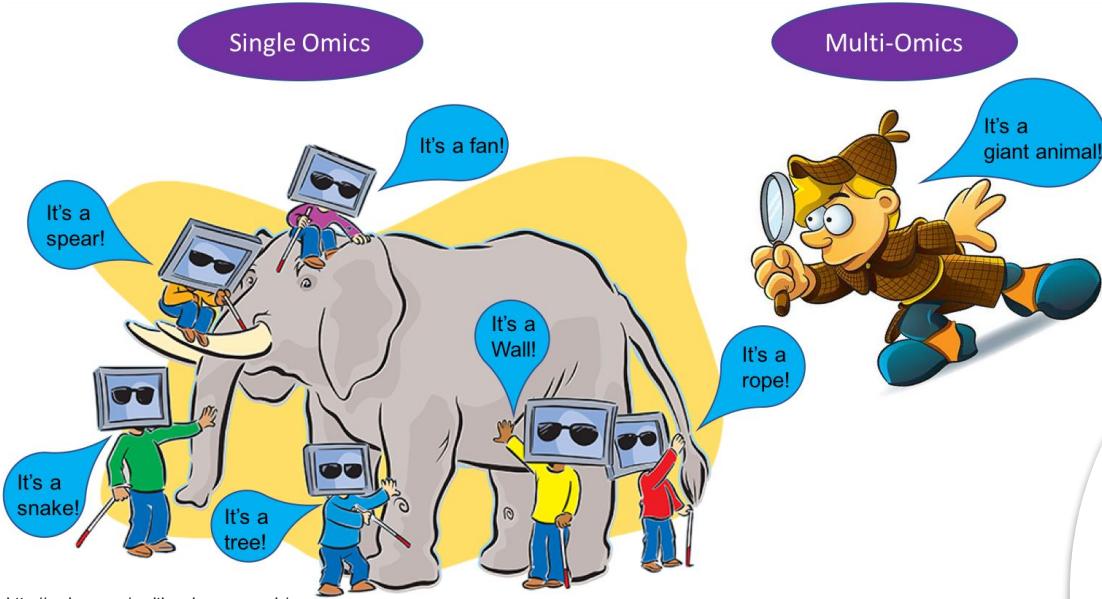
Qu'est-ce que le *multi-omique* ?



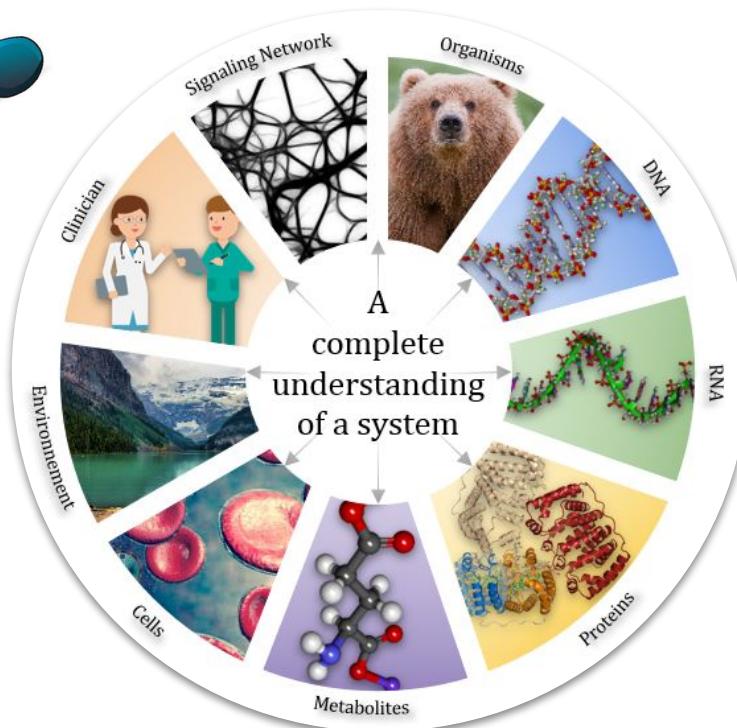
Qu'est-ce que le *multi-omique* ?



Qu'est-ce que le *multi-omique* ?

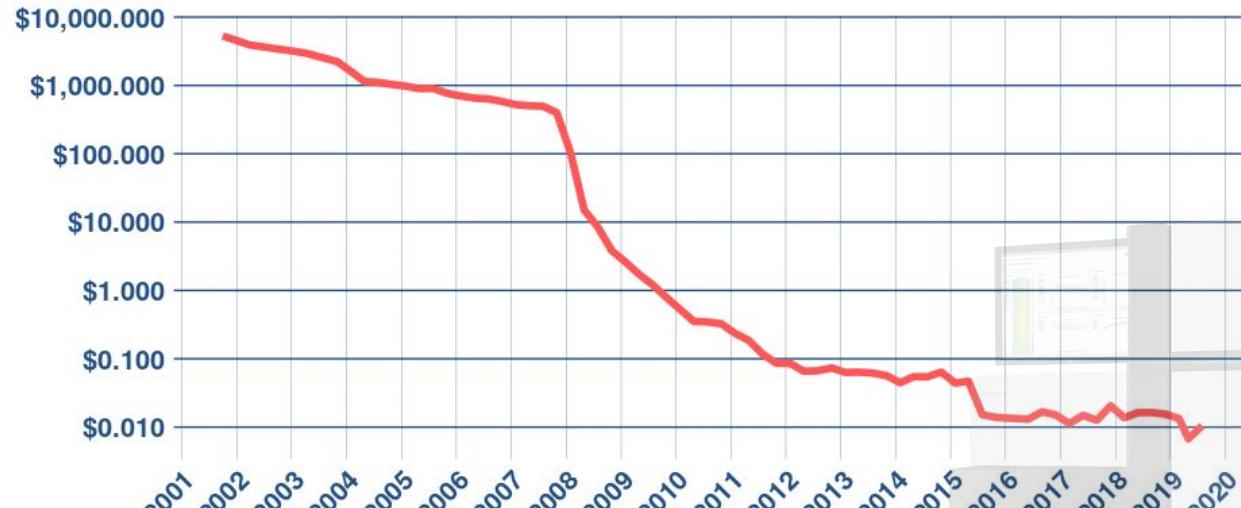


Approche **Réductionniste** (single-omic)
vs
Approche **Holistique** (multi-omic)



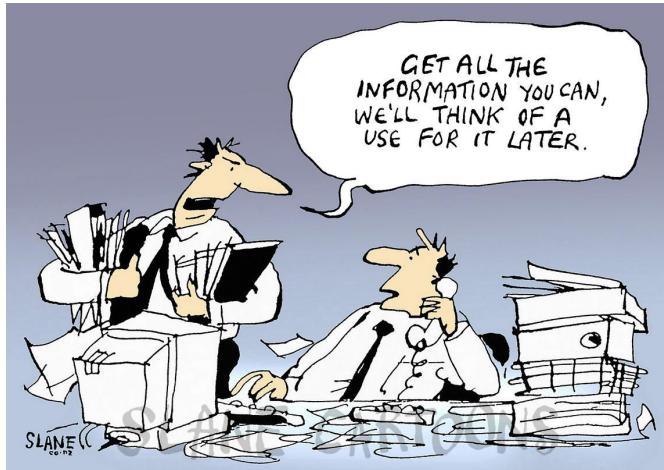
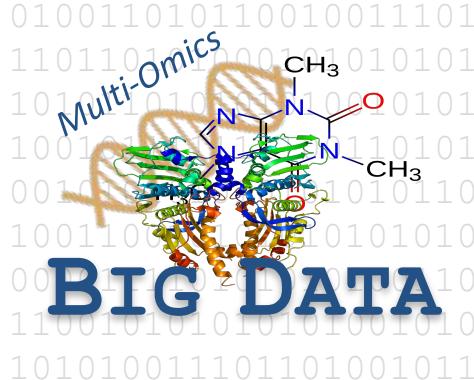
Pourquoi ?

Cost per Raw Megabase of DNA Sequence



genome.gov/sequencingcosts

Difficultés



- Sources multiples et hétérogènes
- Puissance de calcul
- Interprétabilité
- Rester informé des avancées dans chaque domaines

Quelques ressources multi-omique ?

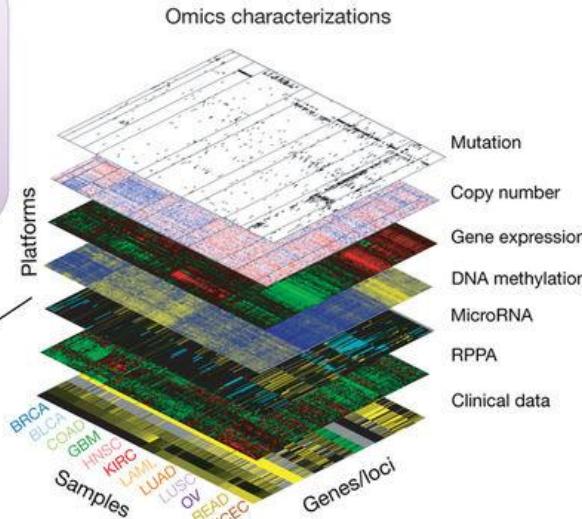
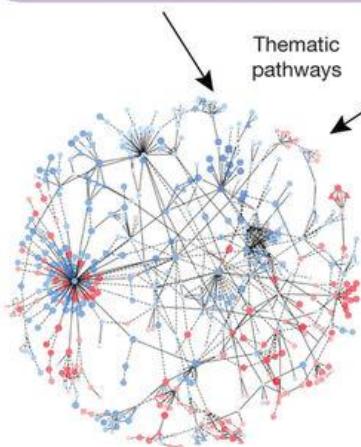
NIH

THE CANCER GENOME ATLAS

National Cancer Institute

National Human Genome Research Institute

12 tumor types



doi:10.1038/ng.2764

TCGA BY THE NUMBERS

TCGA produced over
2.5 PETABYTES
of data



To put this into perspective, 1 petabyte of data is equal to

212,000 DVDs



TCGA data describes
33 DIFFERENT
TUMOR TYPES
...including
10 RARE
CANCERS

...based on paired tumor and normal tissue sets collected from

11,000 PATIENTS

...using
7 DIFFERENT
DATA TYPES



www.cancer.gov/ccg

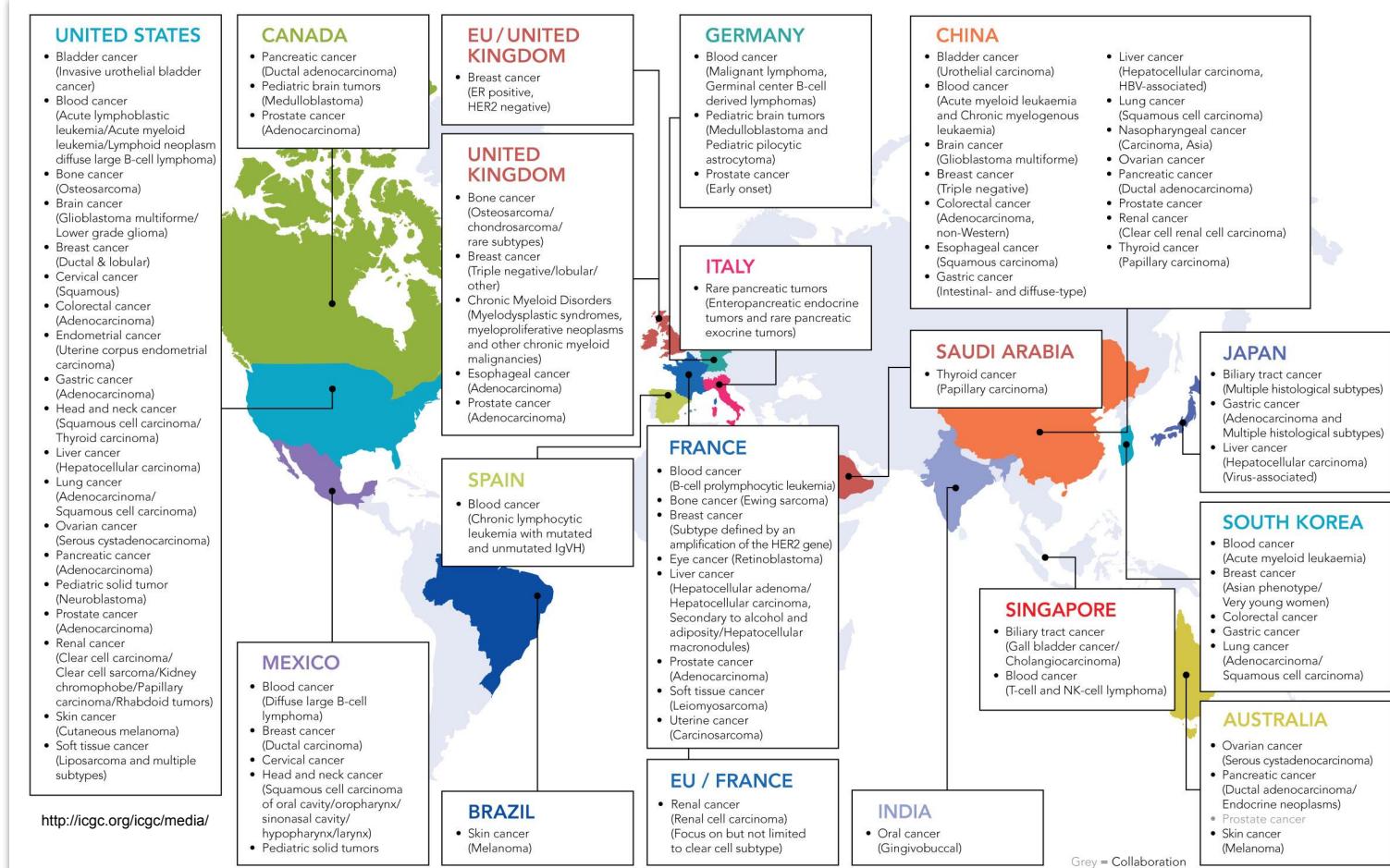
Quelques ressources multi-omique ?



International
Cancer Genome
Consortium



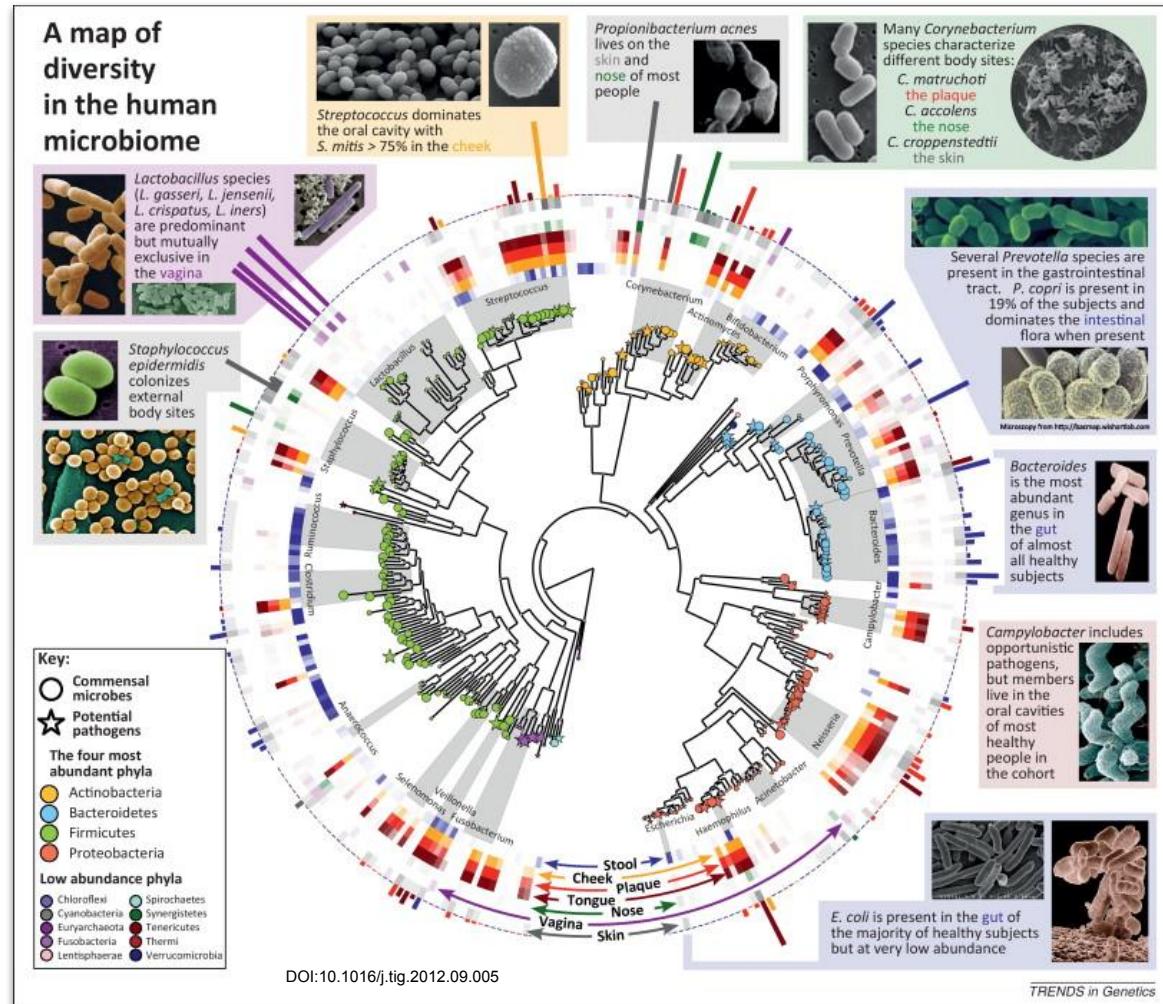
<http://icgc.org/icgc/media/>



Quelques ressources multi-omique ?



<https://hmpdacc.org>





High Throughput
Technology

“Bloc”



Mise en contexte

- omique ? multi-omique ?
- Pourquoi ? Difficultés
- Ressources

Méthodes d'intégration

- Intégration Conceptuelle
- Intégration Statistique
- Intégration à base de modèle

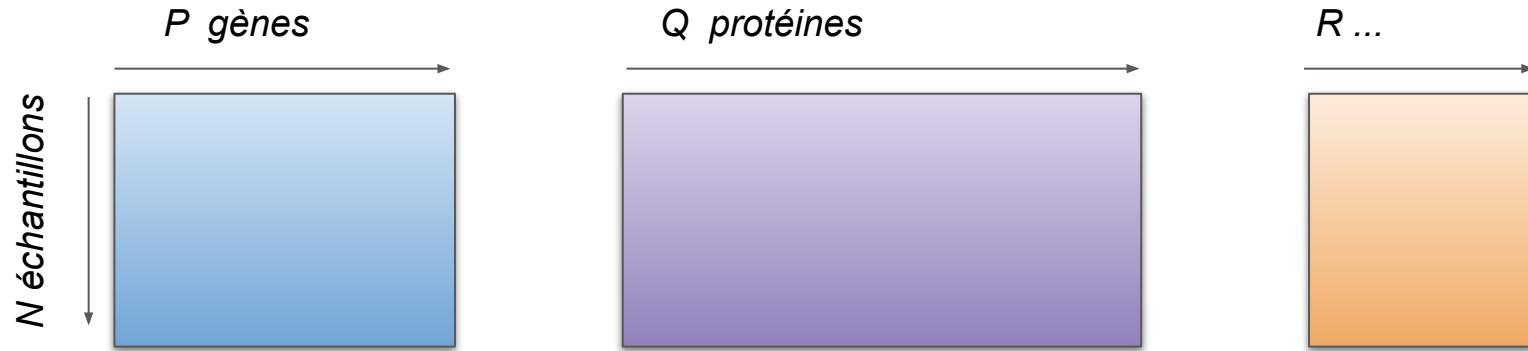
Méthodes Multivariées

- Analyse en Composante Principale
- Autres méthodes multivariées

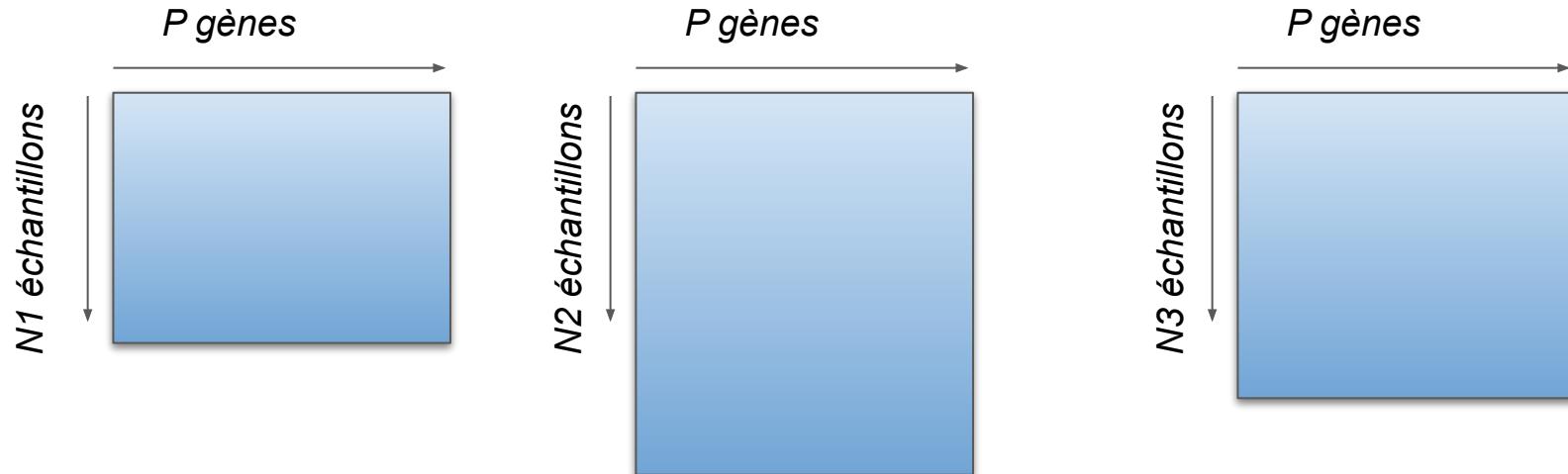
Méthodes à base de réseaux biologiques

- Les réseaux en biologie
- Exploration des réseaux

N-Intégration



P-Intégration



Intégration Conceptuelle



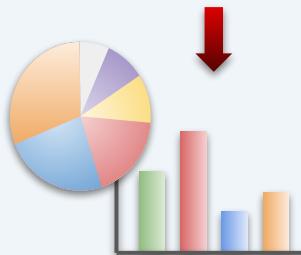
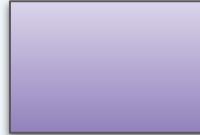
- Complexité +



Données Omiques 1



Données Omiques 2



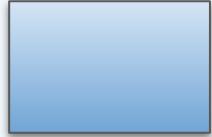
- ❑ Analyse de chaque omique séparément
- ❑ Manque les relations entre les données

Intégration Statistique

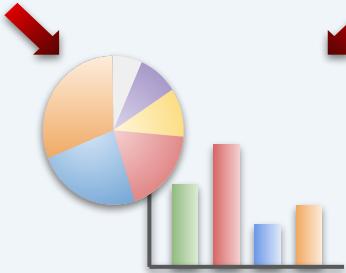
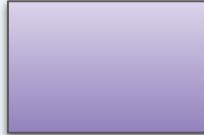
- Complexité +



Données Omiques 1



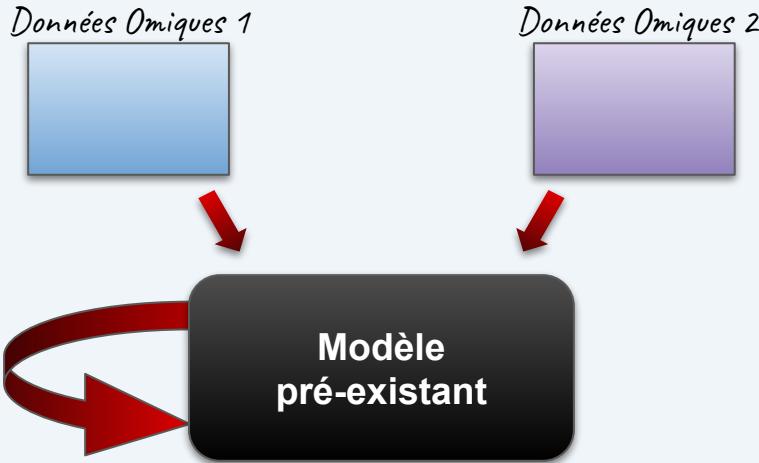
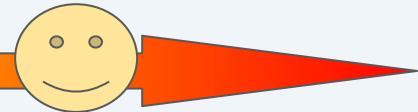
Données Omiques 2



- ❑ Forme d'intégration la plus commune
- ❑ Différentes méthodes :
 - Corrélation
 - Concaténation
 - Multivariée
 - Pathways Biologiques

Intégration à base de Modèles

- Complexité +

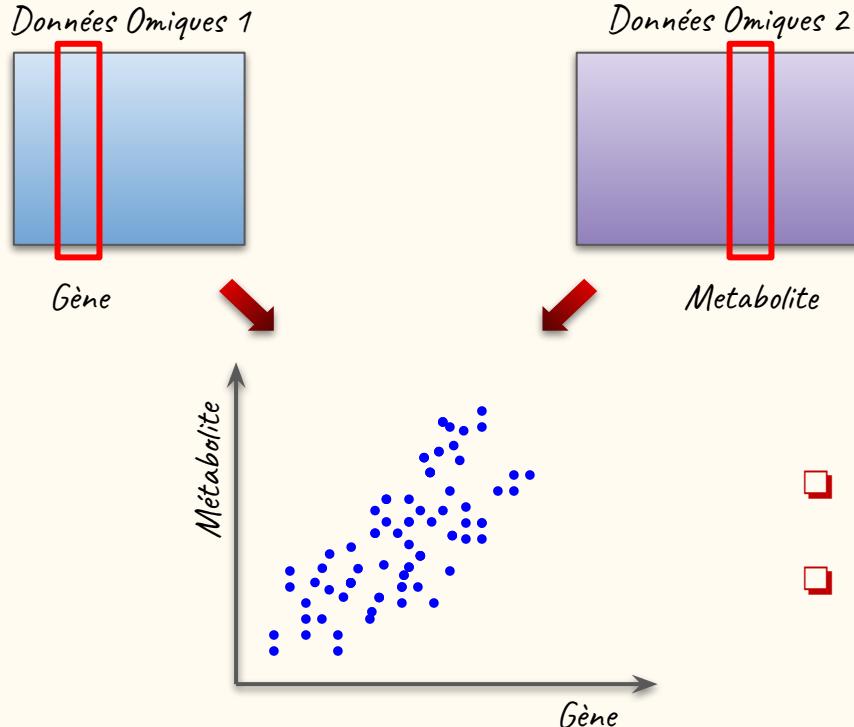


Idéal à atteindre ...

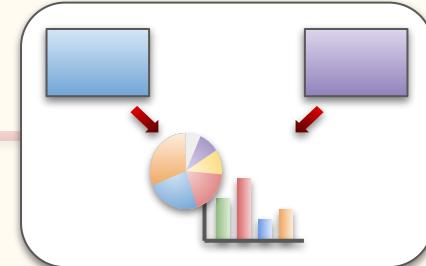
... mais difficilement atteignable aujourd'hui

Intégration Statistique

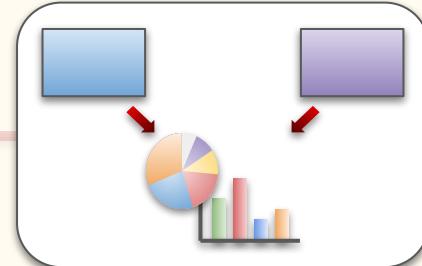
A) Corrélation



- Méthode la plus simple et plus utilisée
- Corrélation de **Pearson** ou **Spearman**



Intégration Statistique

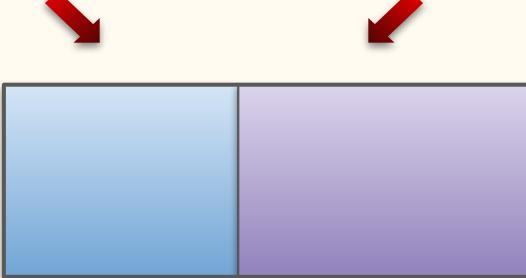


B) Concaténation

Données Omiques 1



Données Omiques 2



- ❑ Concaténation des données en une table
- ❑ **Machine Learning**
- ❑ Mais : poids des données, sources hétérogènes, schéma de valeurs attendues,

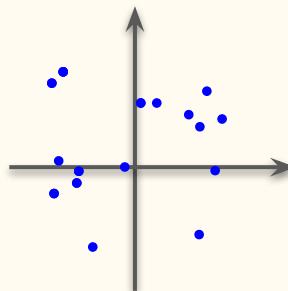
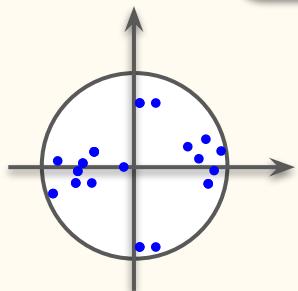
Intégration Statistique

C) Multivarié

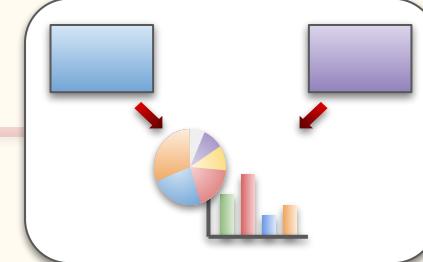
Données Omiques 1



Données Omiques 2



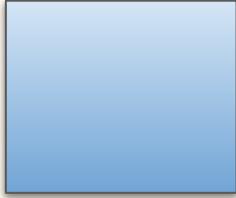
- ❑ Méthodes puissantes et à la mode
- ❑ Technique de PCA, PLS, ...
- ❑ Covariance, correlation



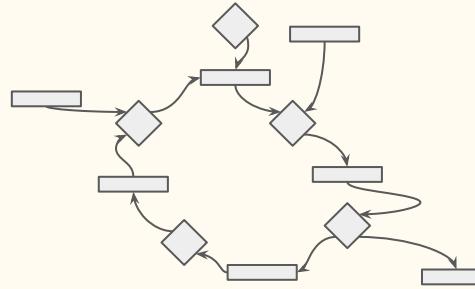
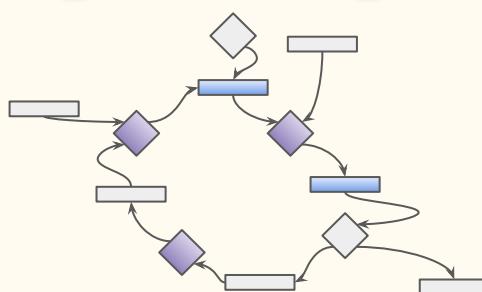
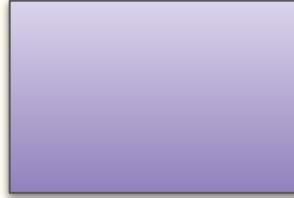
Intégration Statistique

D) Pathways Biologiques

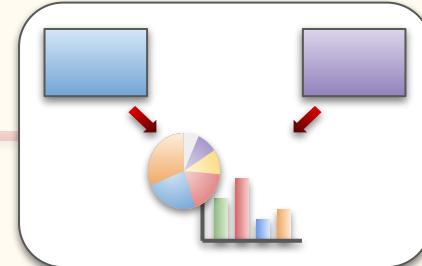
Données Omiques 1



Données Omiques 2



- ❑ “Knowledge driven”
- ❑ Modules fonctionnelles
- ❑ Enrichissement



Mise en contexte

- omique ? multi-omique ?
- Pourquoi ? Difficultés
- Ressources

Méthodes d'intégration

- Intégration Conceptuelle
- Intégration Statistique
- Intégration à base de modèle

Méthodes Multivariées

- Analyse en Composante Principale
- Autres méthodes multivariées

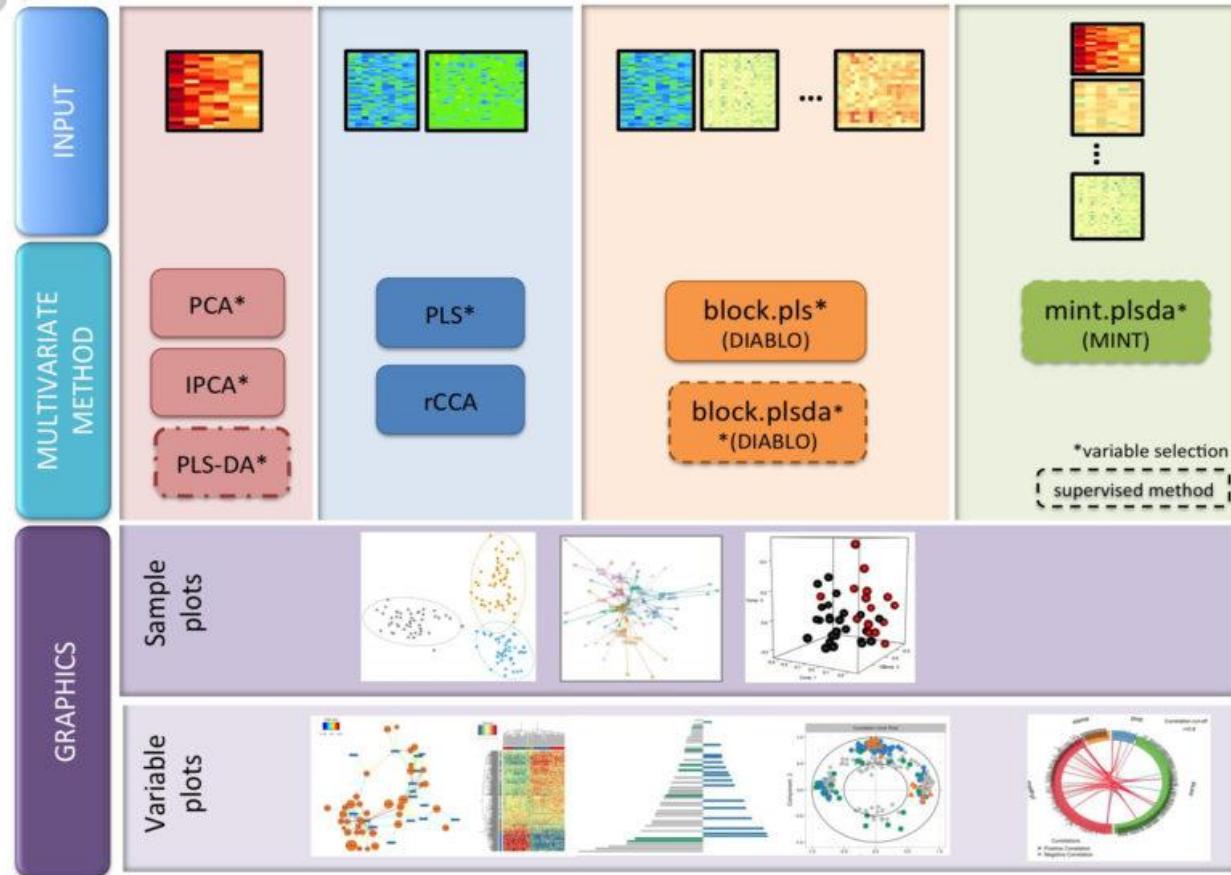
Méthodes à base de réseaux biologiques

- Les réseaux en biologie
- Exploration des réseaux

Choix des méthodes Multivariées



Exploration and
Integration of
Omics datasets



L'Analyse en Composante Principale

Les données

Men

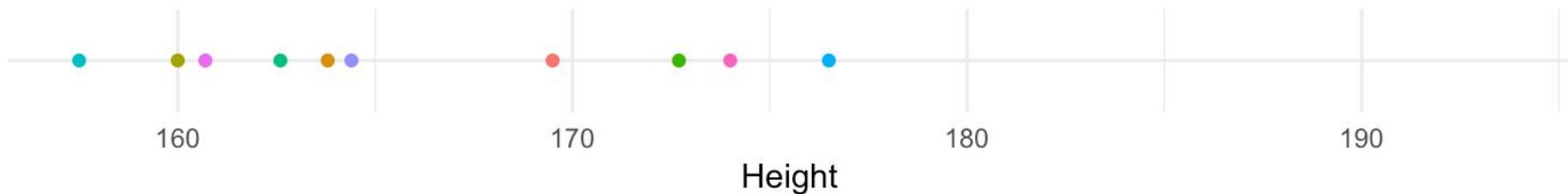
ID	Shoulders	Chest	Waist	Mass	Height
M1	106.2	89.5	71.5	65.6	174.0
M2	110.5	97.0	79.0	71.8	175.3
M3	115.1	97.5	83.2	80.7	193.5
M4	104.5	97.0	77.8	72.6	186.5
M5	107.5	97.5	80.0	78.8	187.2
M6	119.8	99.9	82.5	74.8	181.5
M7	123.5	106.9	82.0	86.4	184.0
M8	120.4	102.5	76.8	78.4	184.5
M9	111.0	91.0	68.5	62.0	175.0
M10	119.5	93.5	77.5	81.6	184.0

Women

ID	Shoulders	Chest	Waist	Mass	Height
W1	105.0	89.0	71.2	67.3	169.5
W2	100.2	94.1	79.6	75.5	160.0
W3	99.1	90.8	77.9	68.2	172.7
W4	107.6	97.0	69.6	61.4	162.6
W5	104.0	95.4	86.0	76.8	157.5
W6	108.4	91.8	69.9	71.8	176.5
W7	99.3	87.3	63.5	55.5	164.4
W8	91.9	78.1	57.9	48.6	160.7
W9	107.1	90.9	72.2	66.4	174.0
W10	100.5	97.1	80.4	67.3	163.8

L'Analyse en Composante Principale

Comment représenter les données ?



L'Analyse en Composante Principale

Indicateur de dispersion

□ La **moyenne** $\bar{x} = \frac{1}{n} \left(\sum_{i=1}^n x_i \right) = \frac{x_1 + x_2 + \cdots + x_n}{n}$

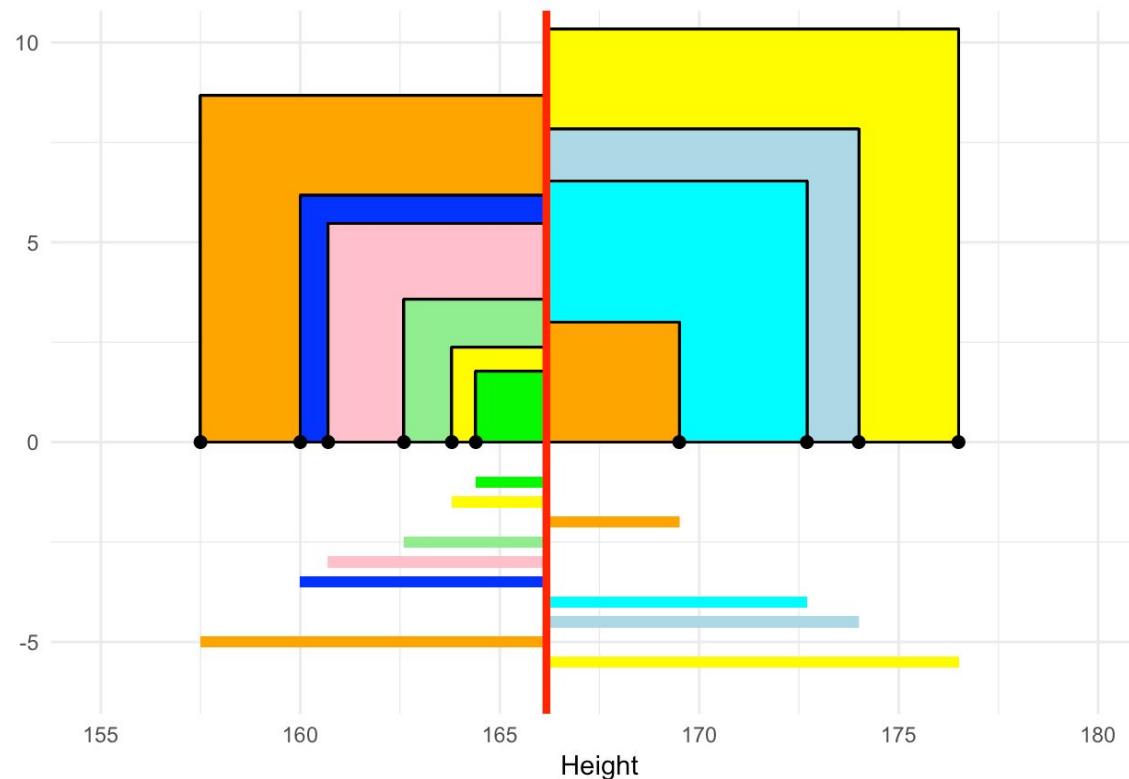
□ L'**étendue** $E = x_{\max} - x_{\min}$

□ La **variance** $V = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$

□ L'**écart-type** $\sigma = \sqrt{V} = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2}$

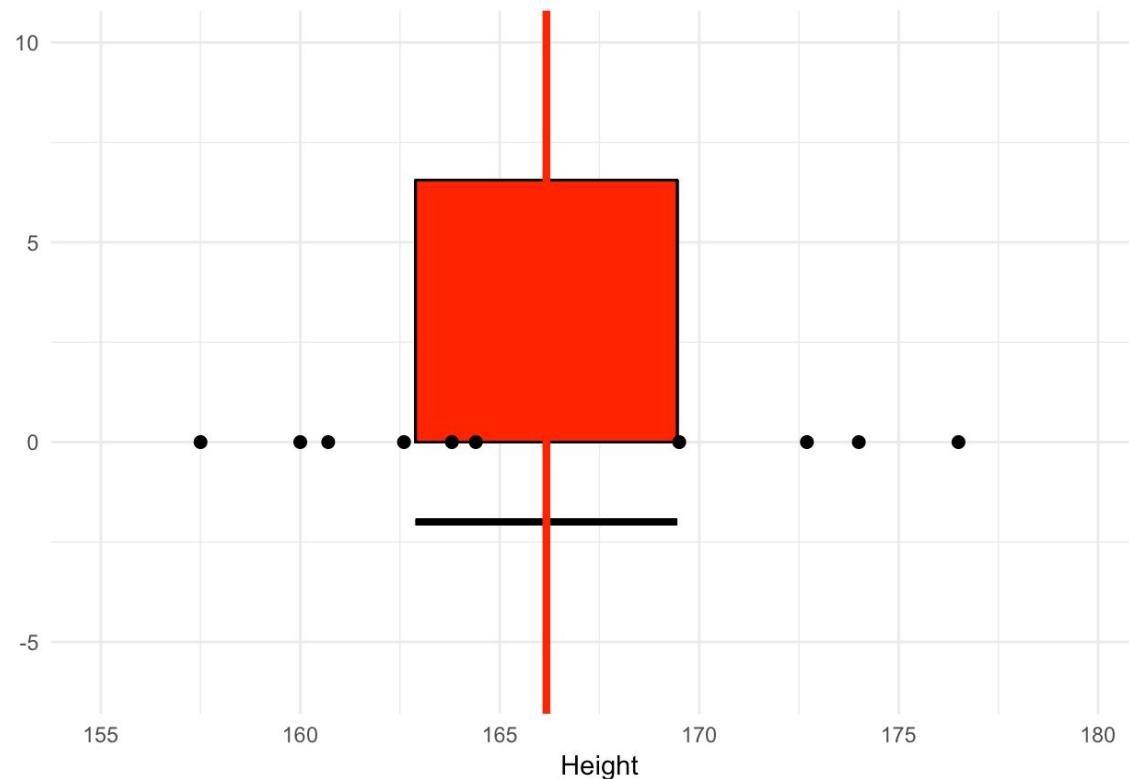
L'Analyse en Composante Principale

Variance et Écart-type



L'Analyse en Composante Principale

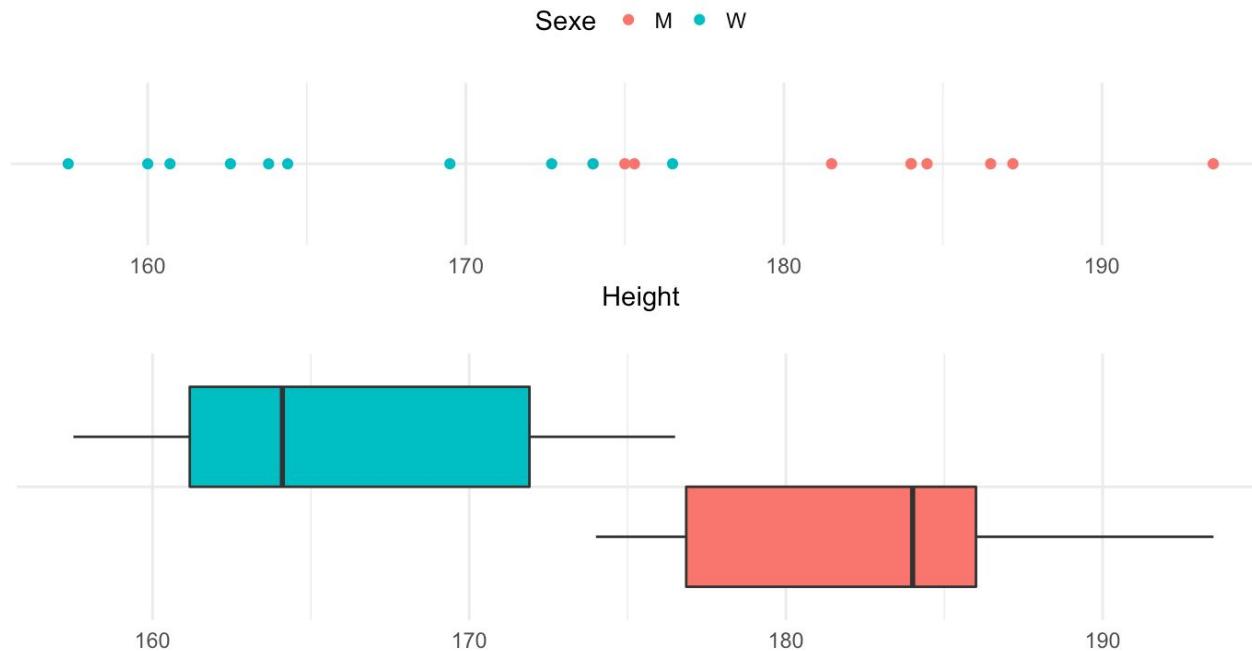
Variance et Écart-type



L'Analyse en Composante Principale

Comment représenter les données ?

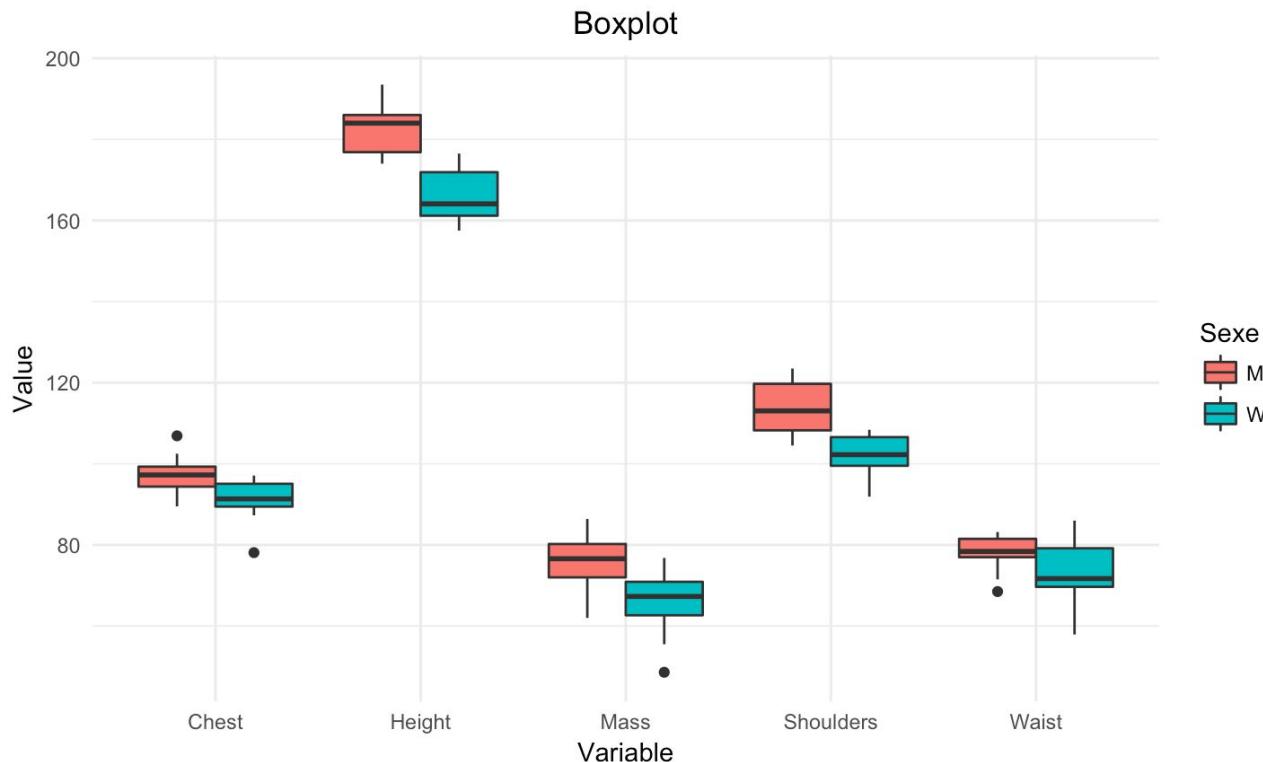
1D



L'Analyse en Composante Principale

Comment représenter les données ?

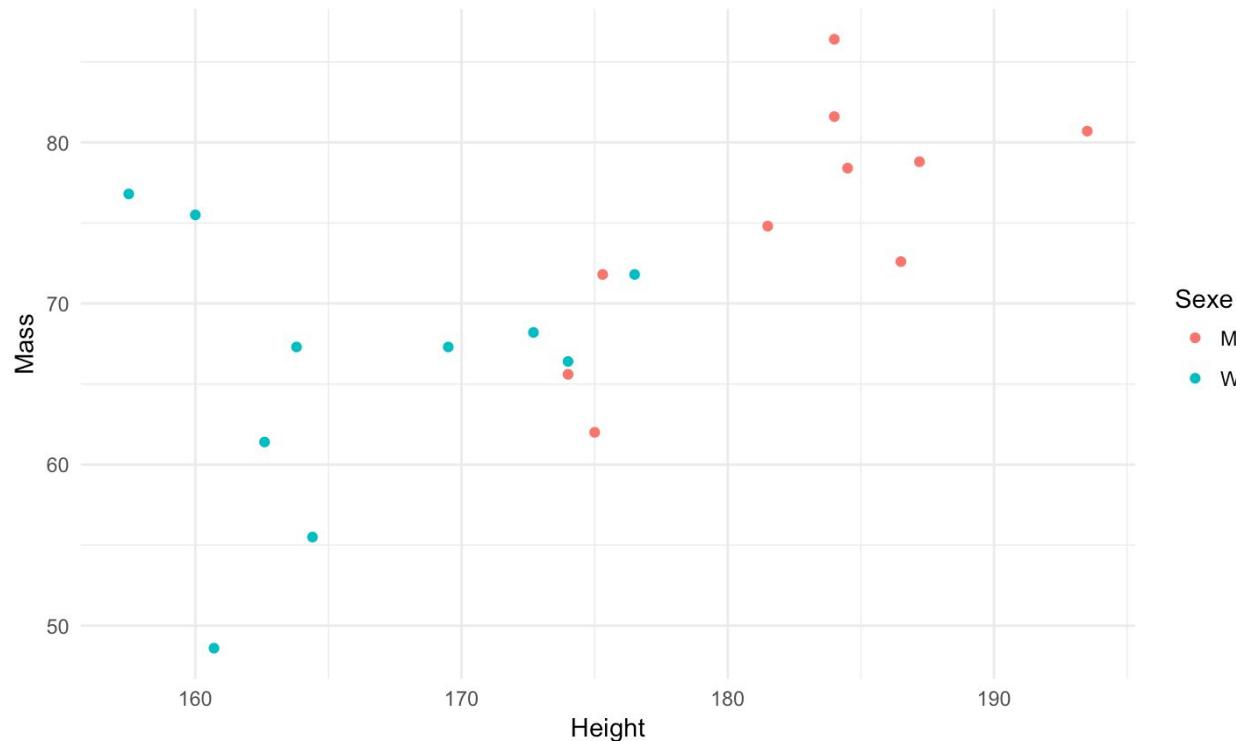
1D



L'Analyse en Composante Principale

Comment représenter les données ?

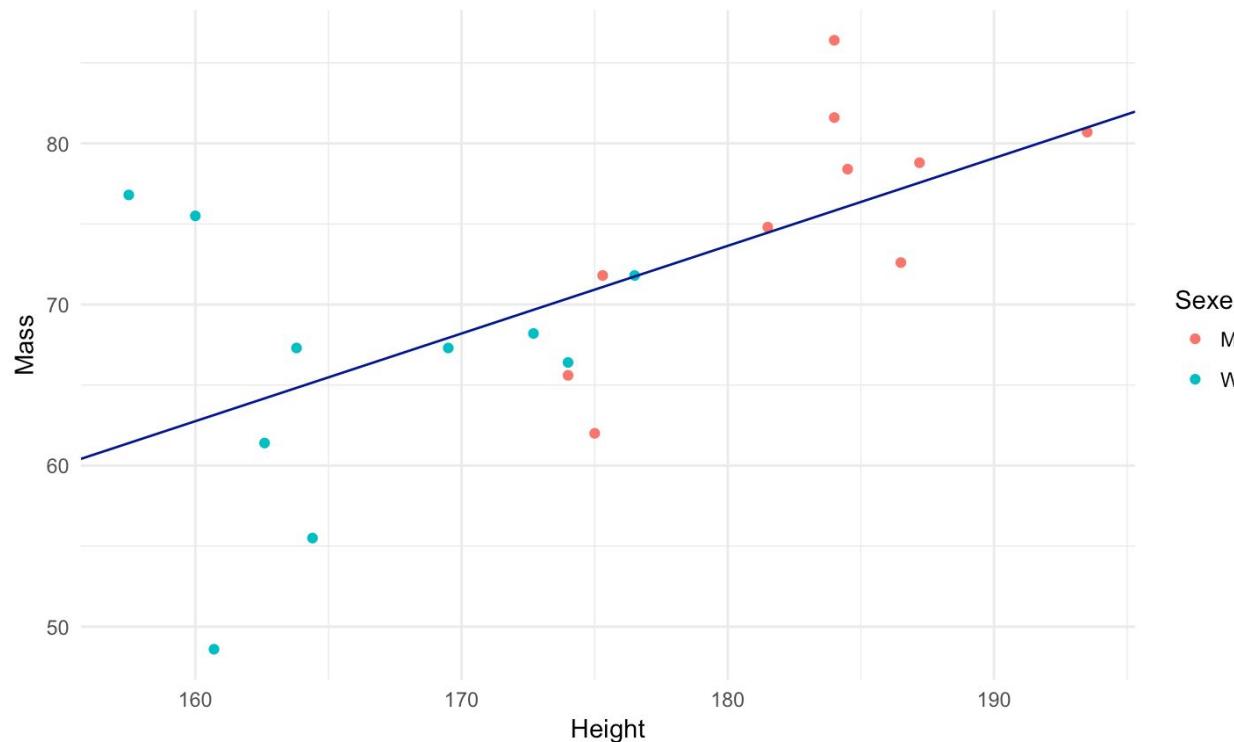
2D



L'Analyse en Composante Principale

Comment représenter les données ?

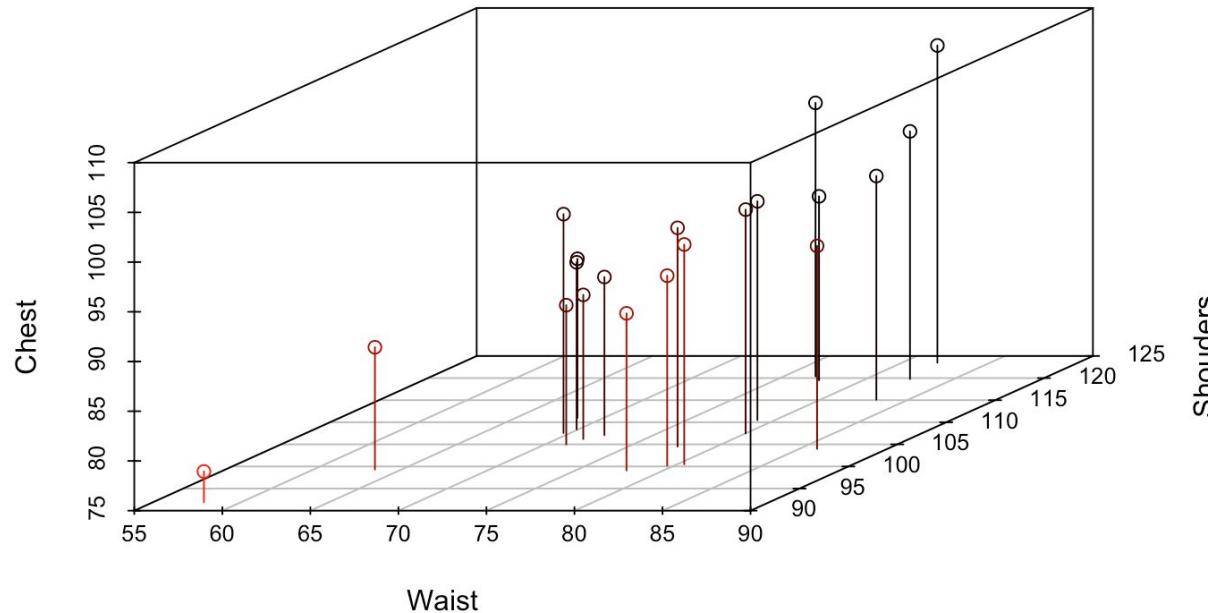
2D



L'Analyse en Composante Principale

Comment représenter les données ?

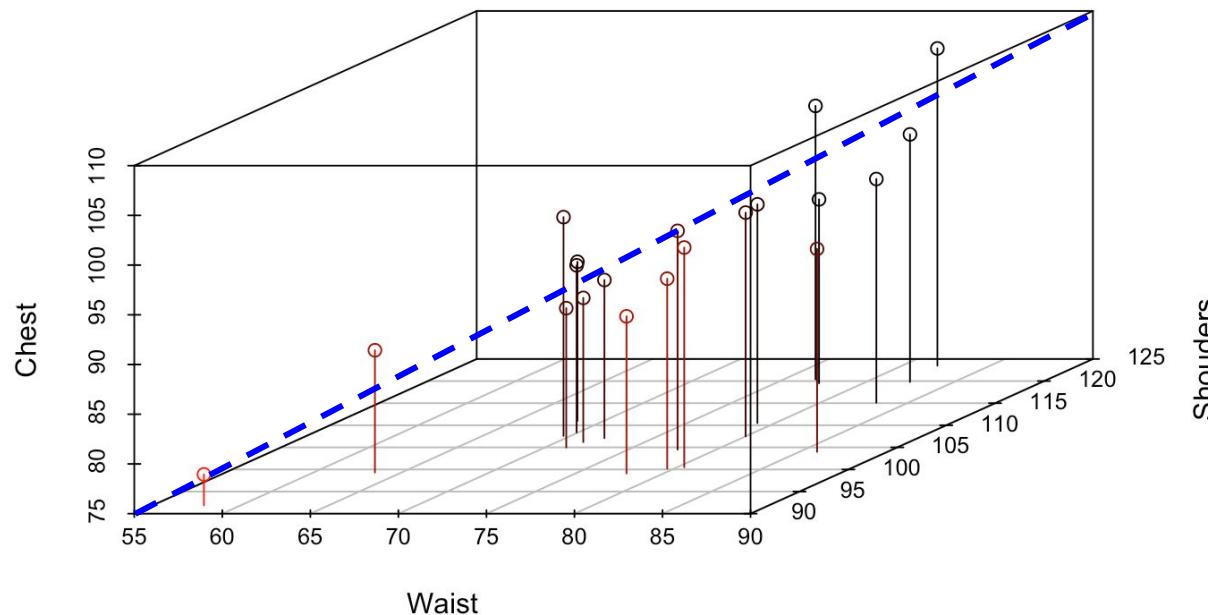
3D



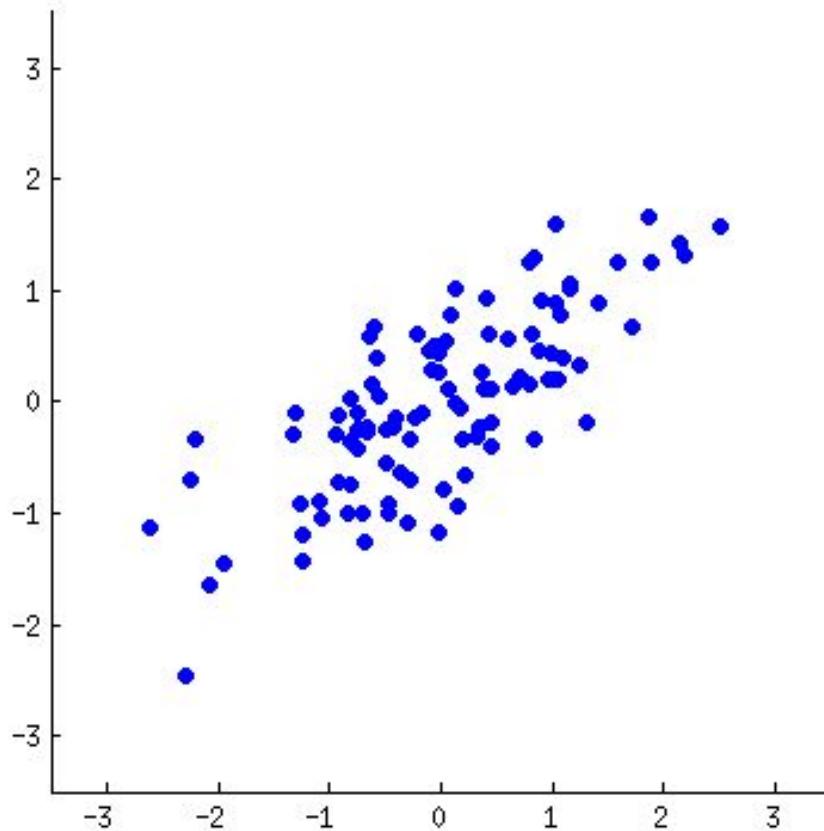
L'Analyse en Composante Principale

Comment représenter les données ?

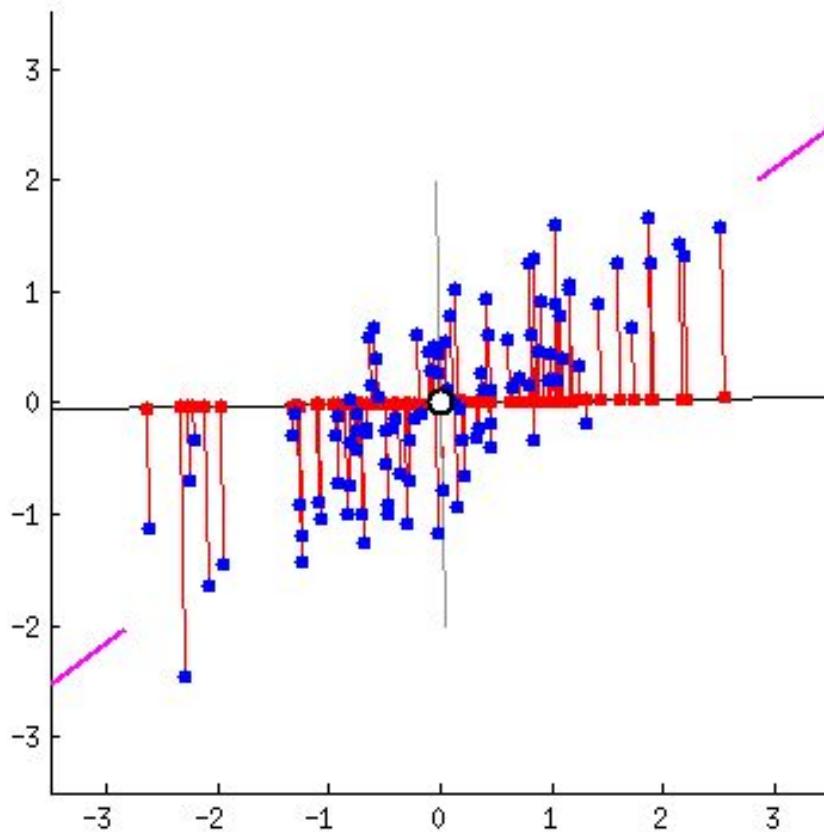
3D



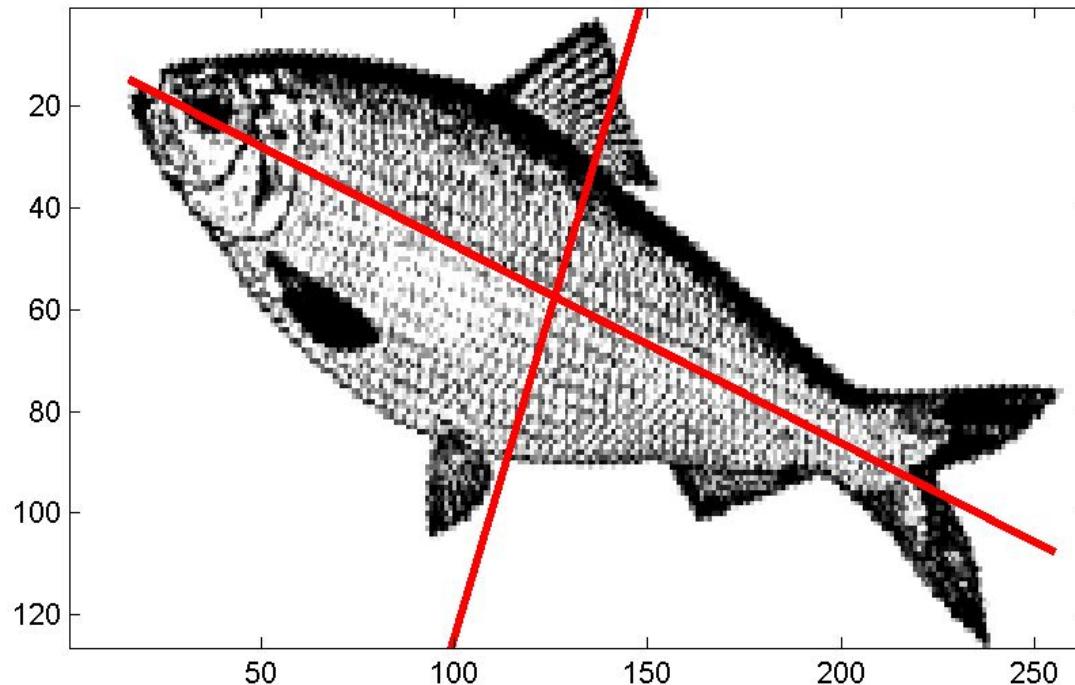
L'Analyse en Composante Principale



L'Analyse en Composante Principale



L'Analyse en Composante Principale



https://upload.wikimedia.org/wikipedia/commons/9/90/PCA_fish.png

L'Analyse en Composante Principale

$X =$

	Gene-A	Gene-B	Gene-C
Ech1	10	5	6
ECh2	42	49	2
...			

components =

	PC1	PC2	PC3
Ech1	6.1	1.2	1.5
ECh2	42.9	8.6	2.5
...			

$$\begin{array}{l} \text{Gene-A} \\ \hline 10 \\ 42 \\ \hline \end{array} + \begin{array}{l} \text{Gene-B} \\ \hline 5 \\ 49 \\ \hline \end{array} + \begin{array}{l} \text{Gene-C} \\ \hline 6 \\ 2 \\ \hline \end{array} = \begin{array}{l} \text{Linear Comb.} \\ \hline 6.1 \\ 42.9 \\ \hline \end{array}$$

$0.2 \times (a1)$

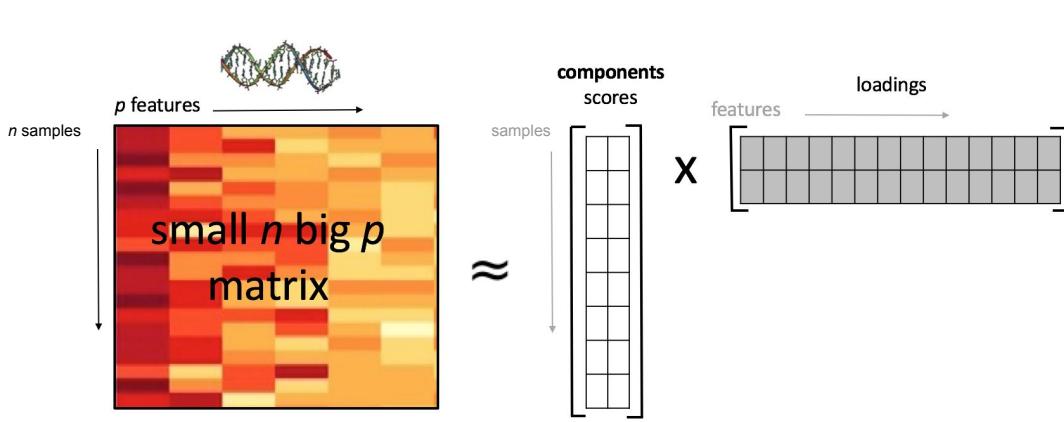
$0.7 \times (a2)$

$0.1 \times (a3)$

$a = \text{loading vectors}$

$= PC$

L'Analyse en Composante Principale



$$\arg \max_{\|a^h\|=1} \text{var}(Xa^h)$$

Solved with SVD:

$$X = U\Delta A^T$$

Singular vectors:

- $T = U\Delta$, T contains the PCs t^h
- A contains the loading vectors a^h

Singular values:

- Δ diagonal matrix with $\sqrt{\delta_h}$

PCA dimension: $h = 1..H$

X
($n \times p$)

is decomposed into:

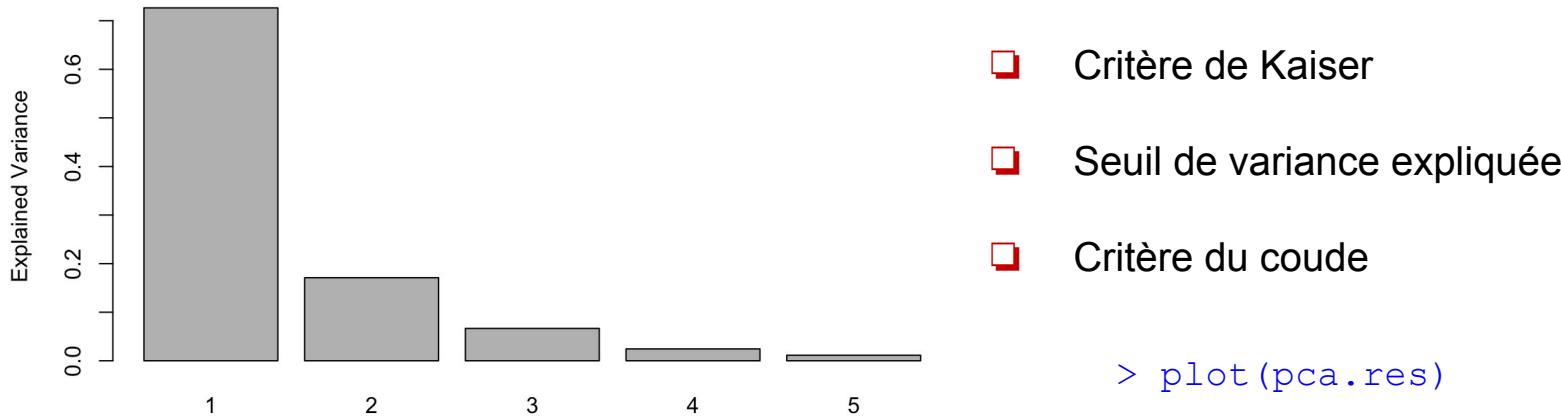
- principal components PC
- loading vectors associated to each PC
=> such that $\text{var}(\text{PC})$ is max.

The variance of the first principal component t^1 is the largest ($= \delta_1$).

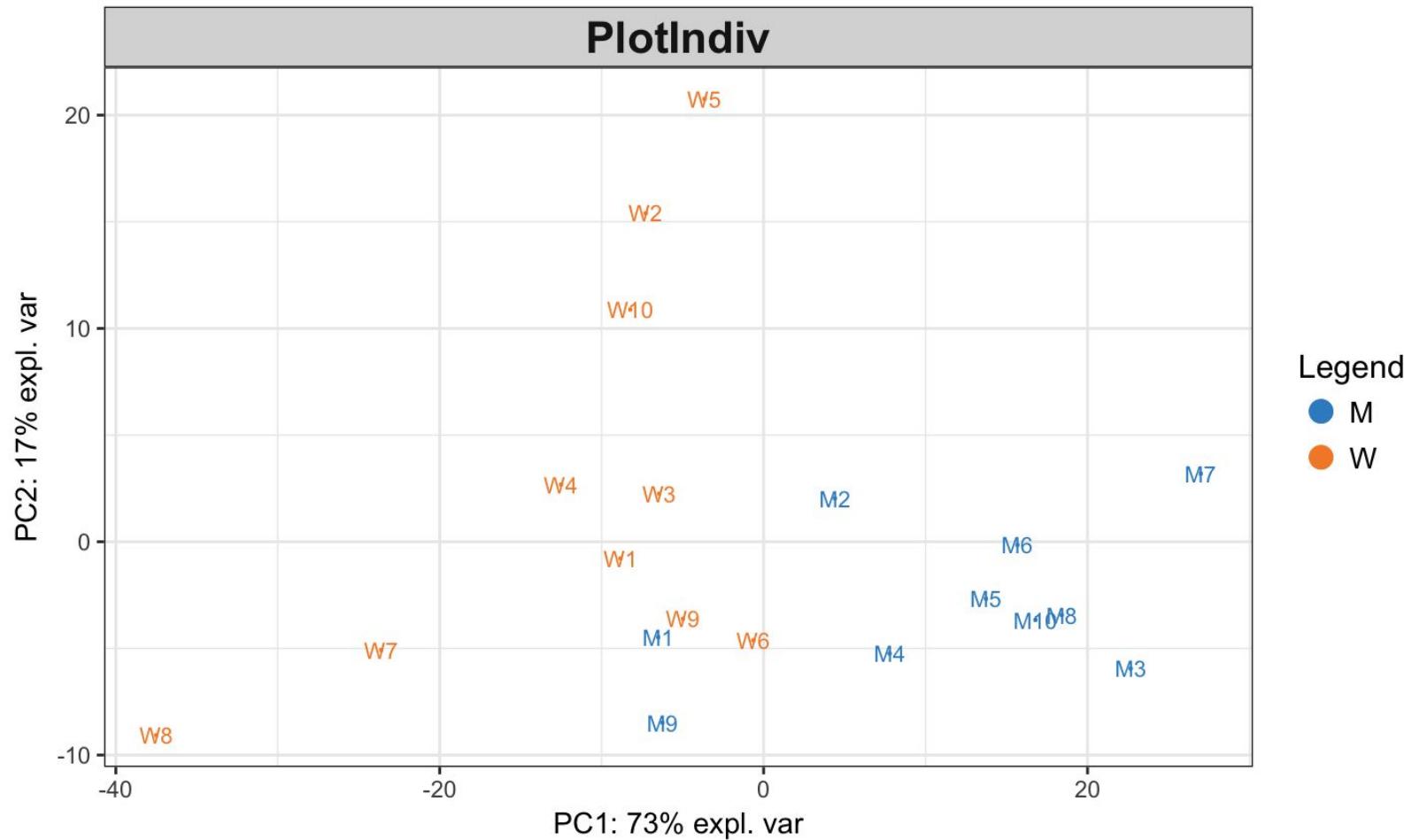
The eigenvalues δ_h decrease and correspond to the explained variance per component.

Exemple d'ACP - Nombre de composantes

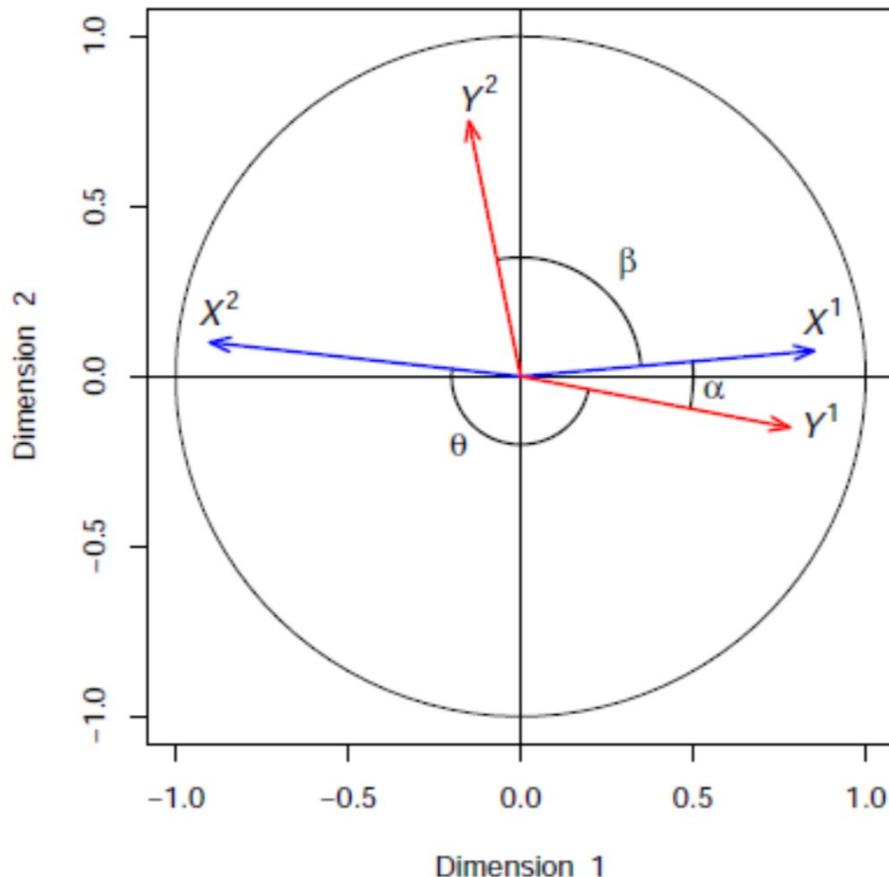
```
## Eigenvalues for the first 5 principal components, see object$sdev^2:  
##          PC1         PC2         PC3         PC4         PC5  
## 255.655833 60.183388 23.484105 8.607990 4.010947  
##  
## Proportion of explained variance for the first 5 principal components, see object$explained_variance:  
##          PC1         PC2         PC3         PC4         PC5  
## 0.72641413 0.17100358 0.06672715 0.02445853 0.01139661  
##  
## Cumulative proportion explained variance for the first 5 principal components, see object$cum.var:  
##          PC1         PC2         PC3         PC4         PC5  
## 0.7264141 0.8974177 0.9641449 0.9886034 1.0000000
```



Exemple d'ACP - graphe des individus

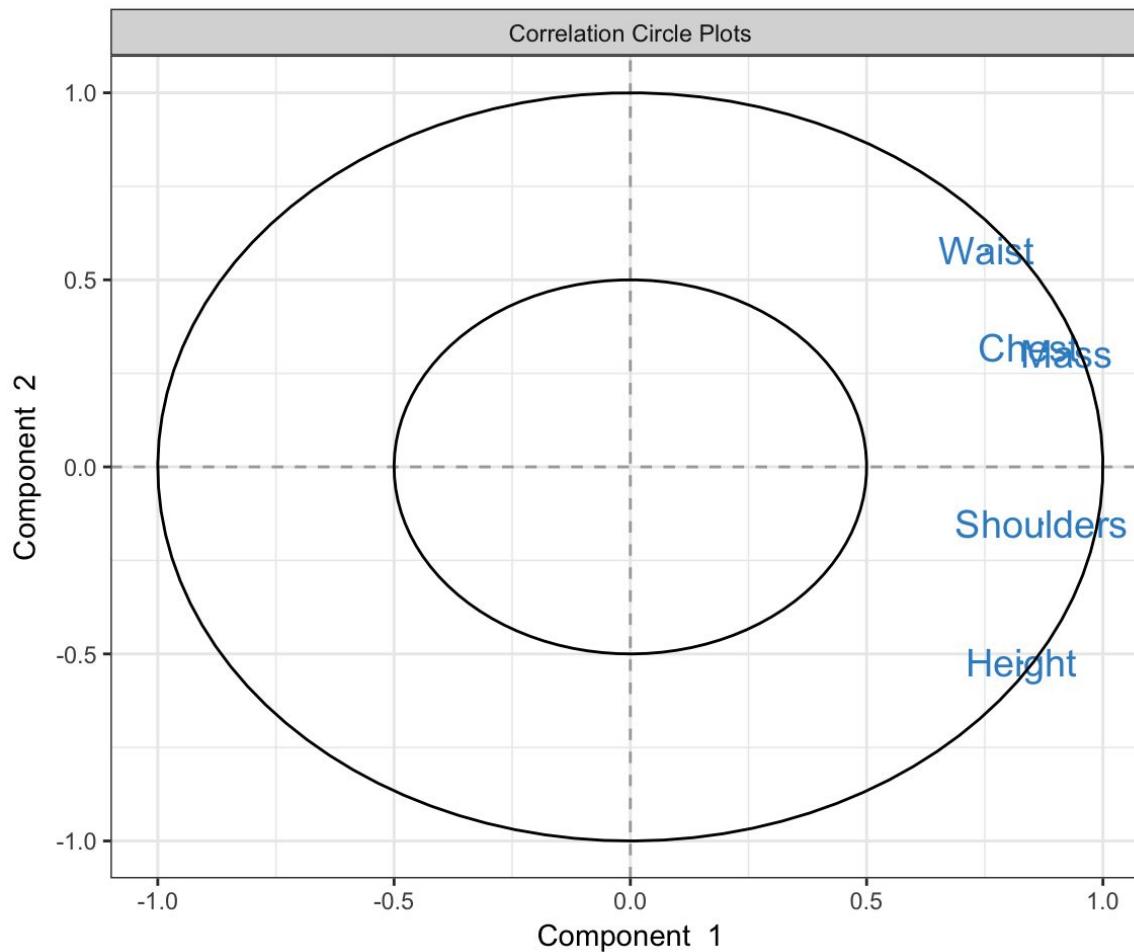


Exemple d'ACP - graphe des variables

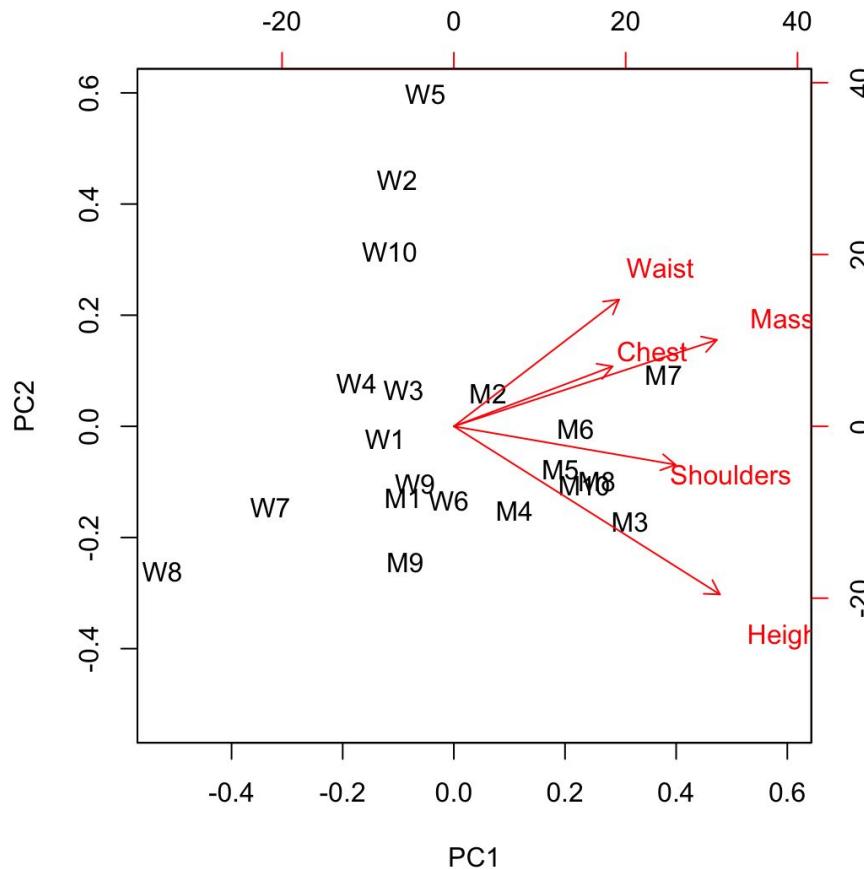


- ❑ Hypersphères des corrélations
- ❑ Corrélation entre 1 variables :
 - ❑ positive
 - ❑ négative
 - ❑ nulle
- ❑ Contribution aux axes

Exemple d'ACP - graphe des variables



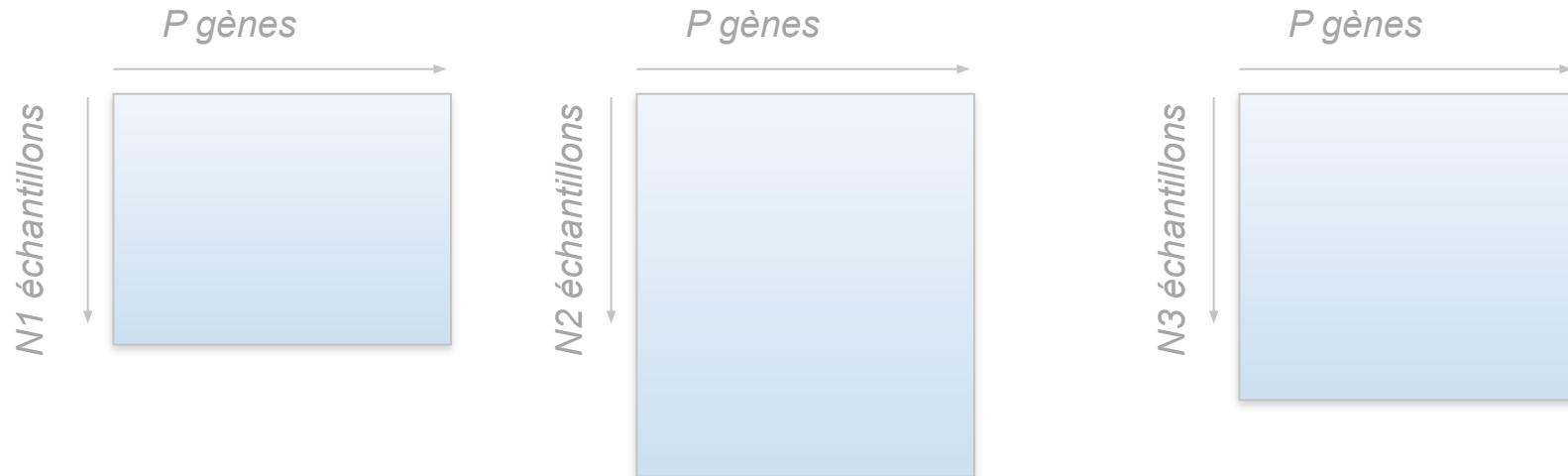
Exemple d'ACP - biplot



N-Intégration



P-Intégration

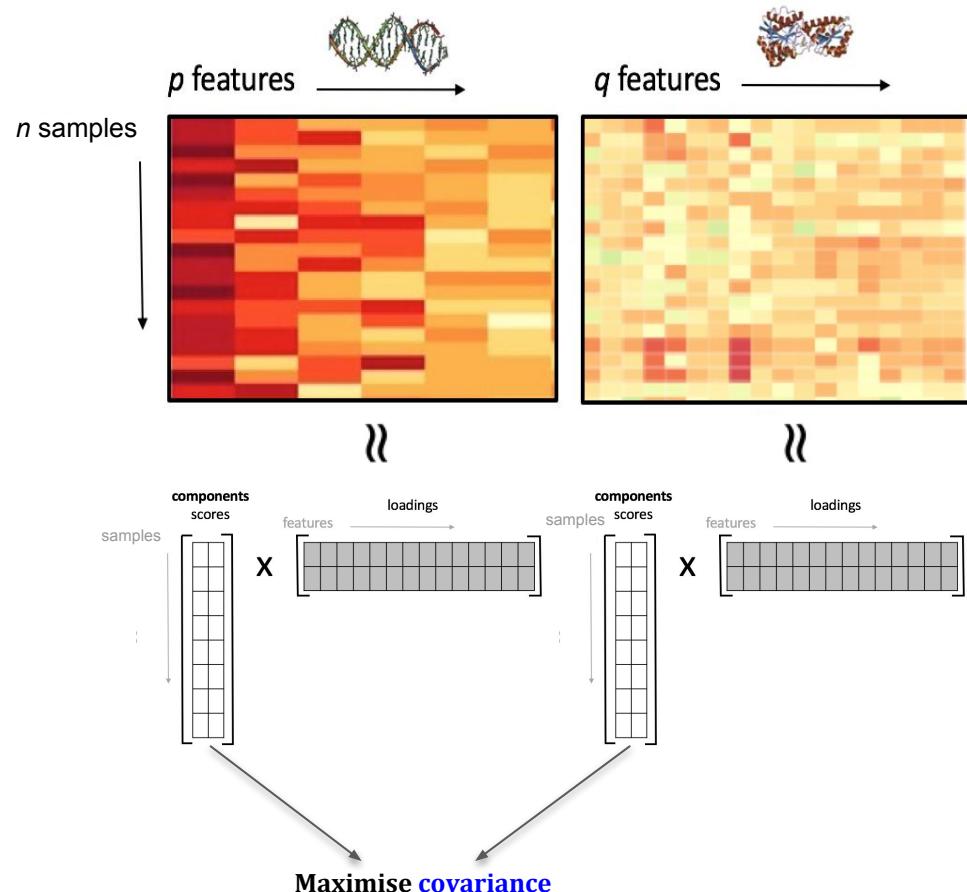


PLS: Projection on Latent Structures

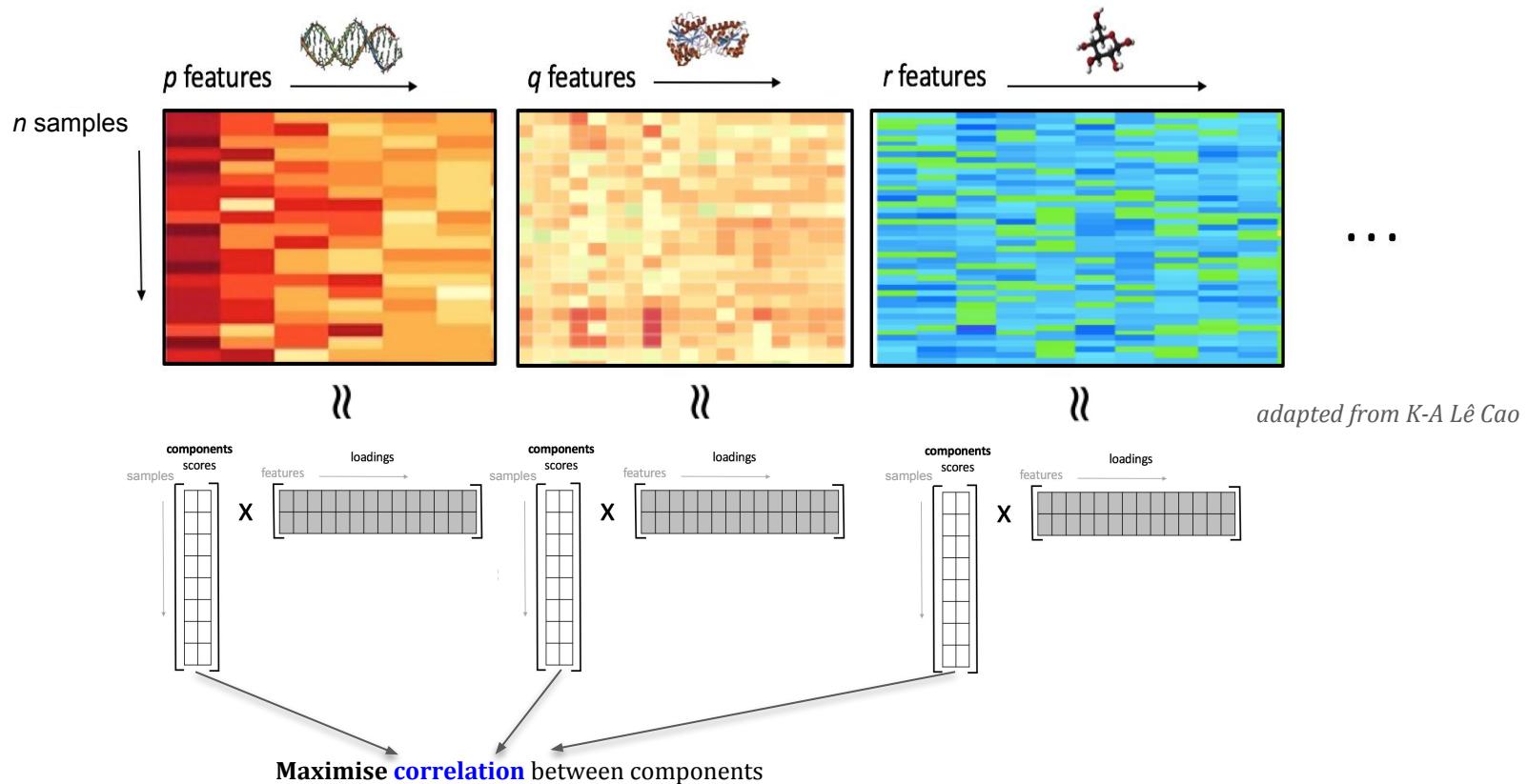
PLS maximises the **covariance** between 2 sets of data

$$\arg \max_{\|a^h\|=1, \|b^h\|=1} \text{cov}(Xa^h, Yb^h)$$

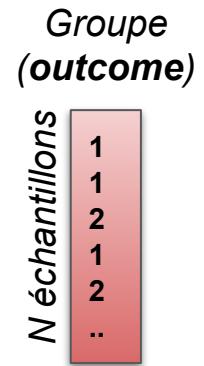
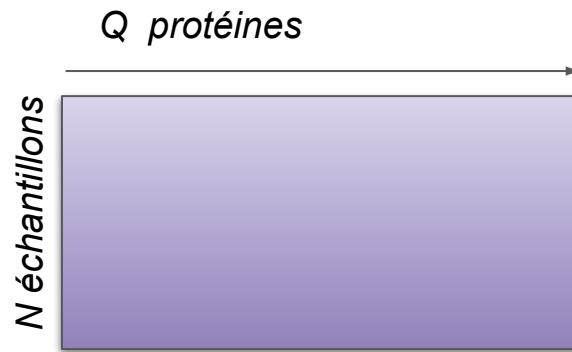
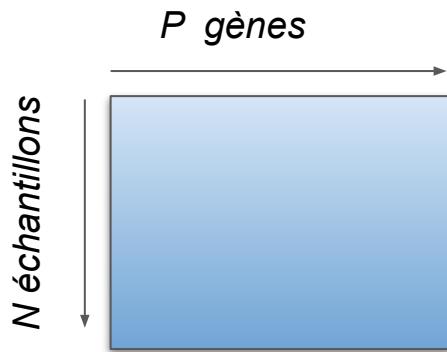
Loading vectors obtained from $\text{svd}(X_1^T X_2)$



Multi-block PLS



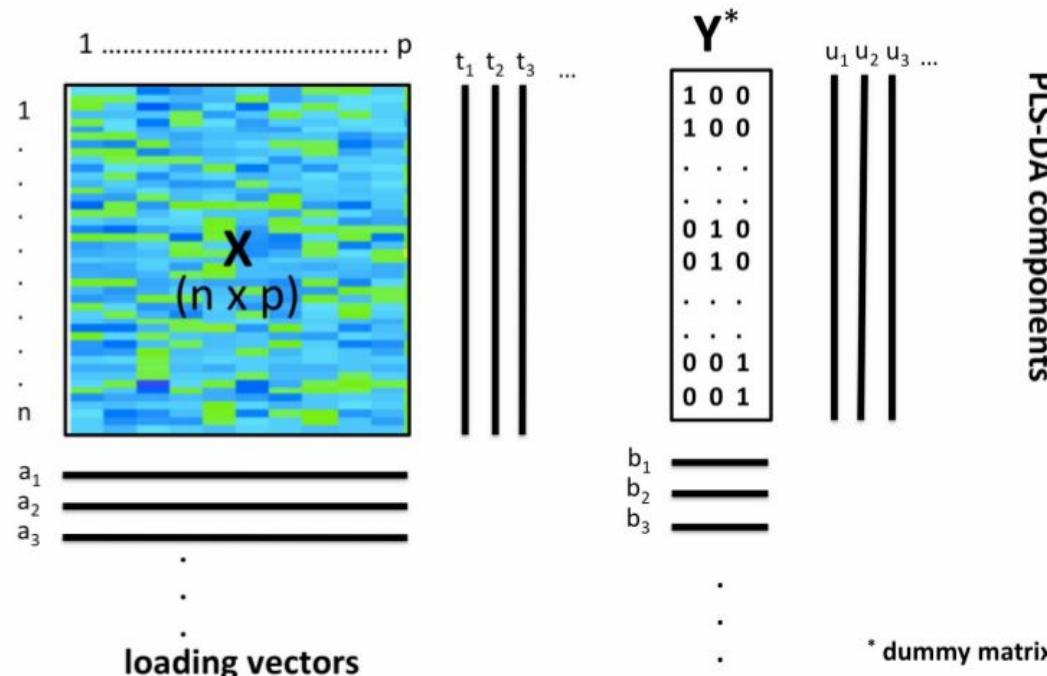
Unsupervised



Supervised

PLS - Discriminant Analysis (PLS-DA)

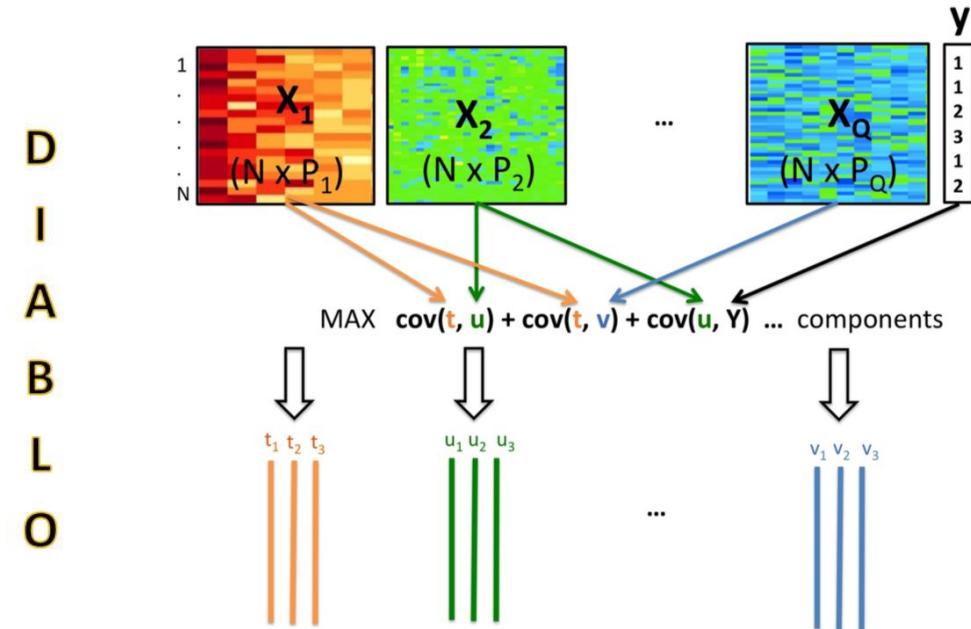
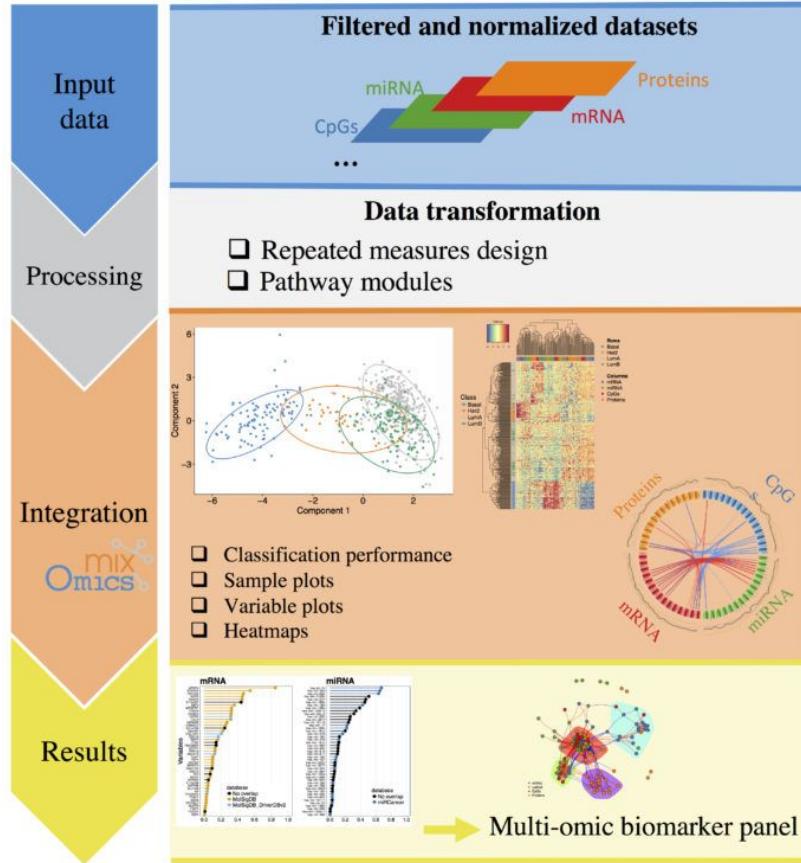
Seek for the PLSDA components from X that best explain the outcome Y^* such that $\text{cov}(t_h, u_h)$ is max.



From K-A Lê Cao

DIABLO

Data Integration Analysis for Biomarker discovery using Latent Variables approaches for Omics studies

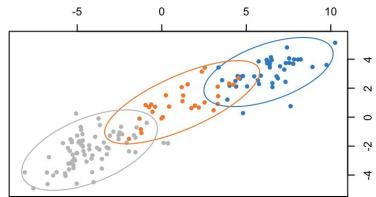


From K-A Lê Cao

DIABLO

mRNA

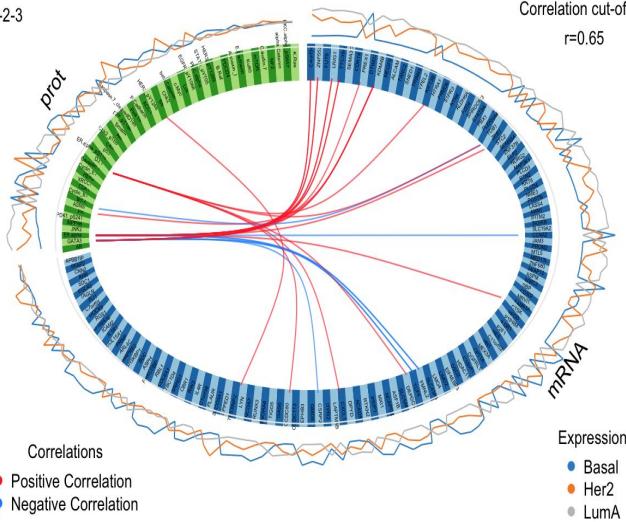
0.93



● Basal ● Her2 ● LumA

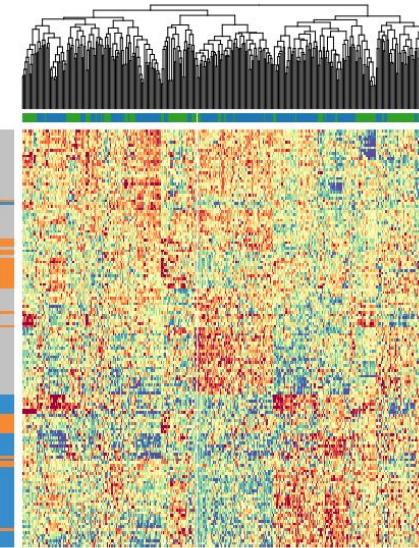
Comp 1-2-3

Correlation cut-off
 $r=0.65$



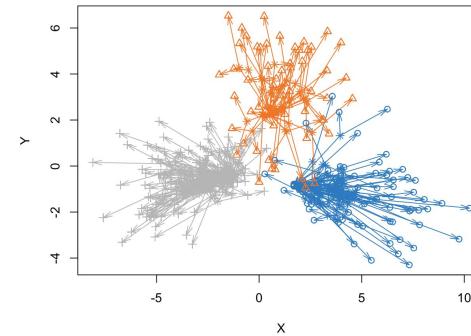
Correlations
● Positive Correlation
● Negative Correlation

Color key
-2 -1 0 1 2

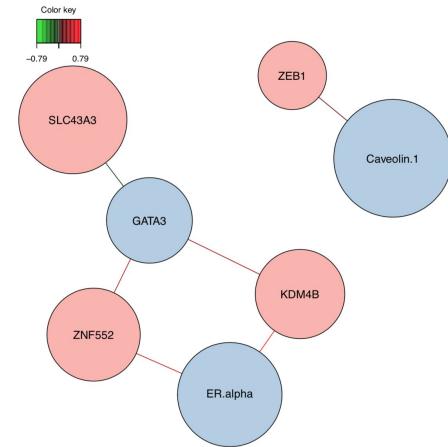


Rows
● Basal
● Her2
● LumA

Columns
● mRNA
● prot

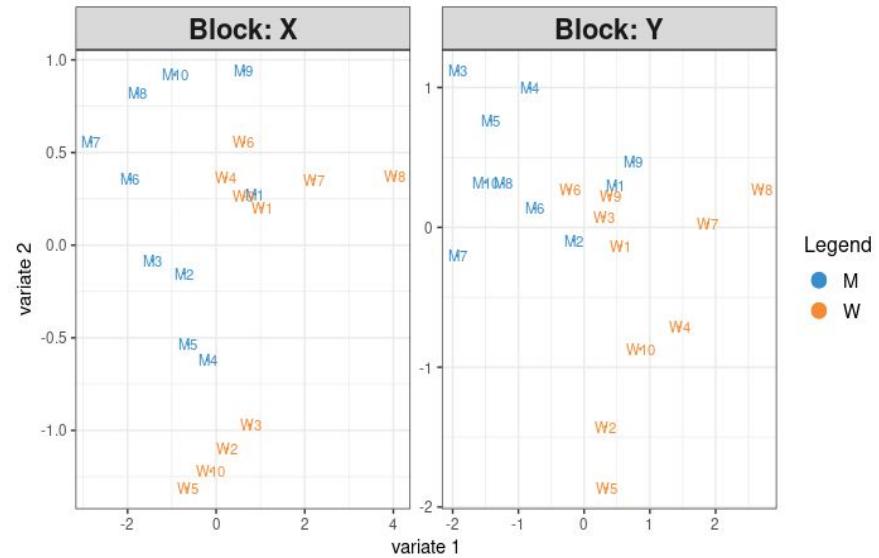
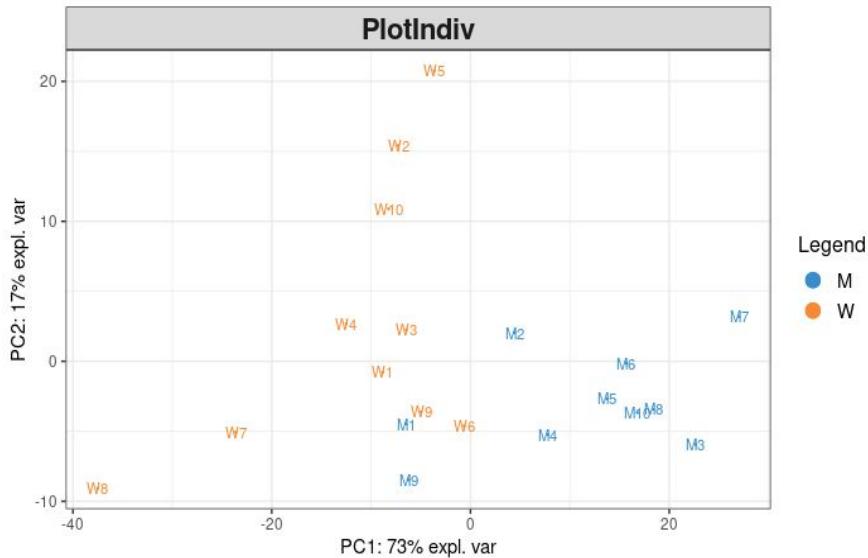


Legend
○ Basal
△ Her2
× LumA



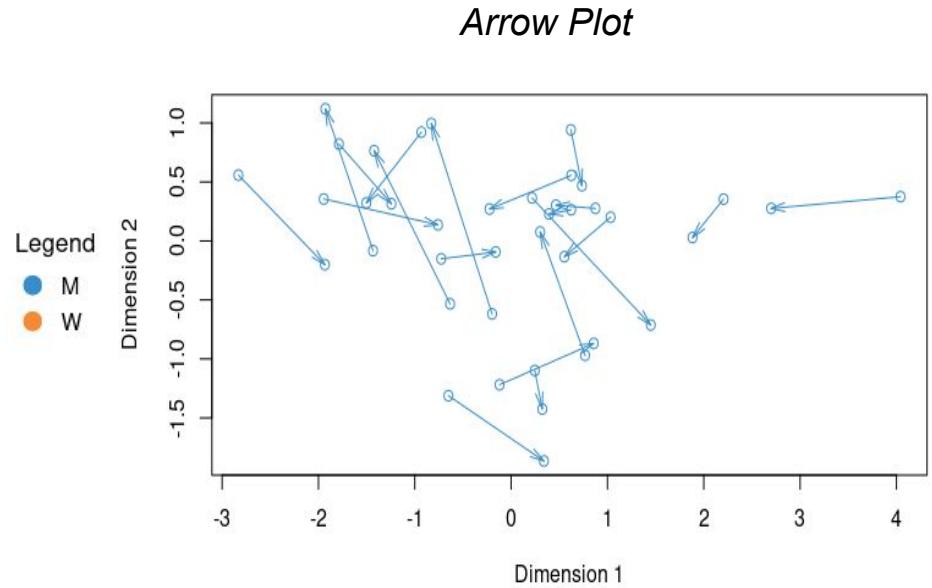
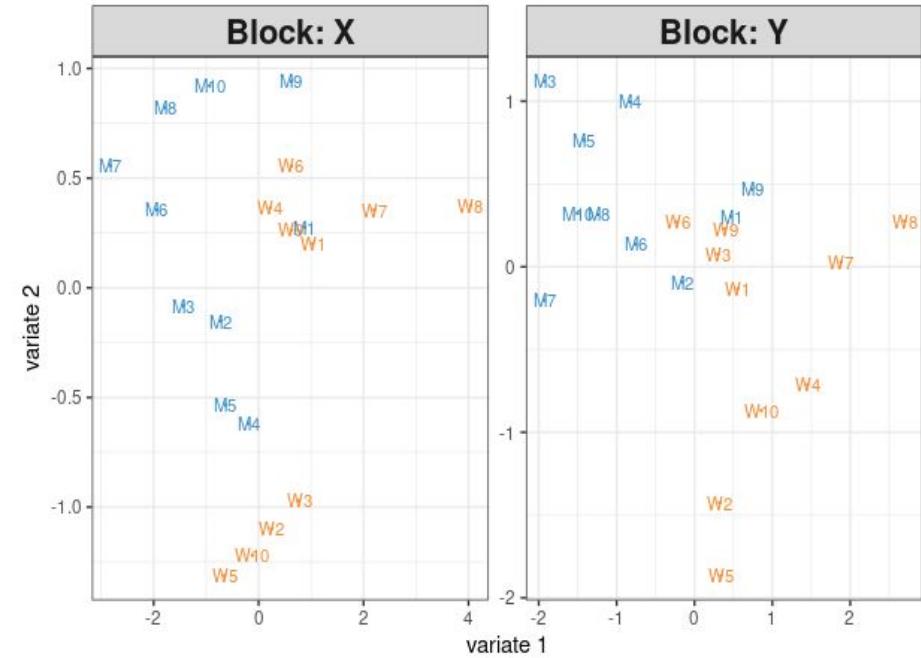
mixOmics Graphical Output

1. Samples plots



- Projection des échantillons (samples) dans l'espace produit par les nouvelles variables instrumentales (components, latent structures, ...)
- similarité entre les échantillons (clustering)
- autres méthodes pour représenter les échantillons dans mixOmics: [plotArrow](#) (paired: X -> Y)

> [plotIndiv\(...\)](#)

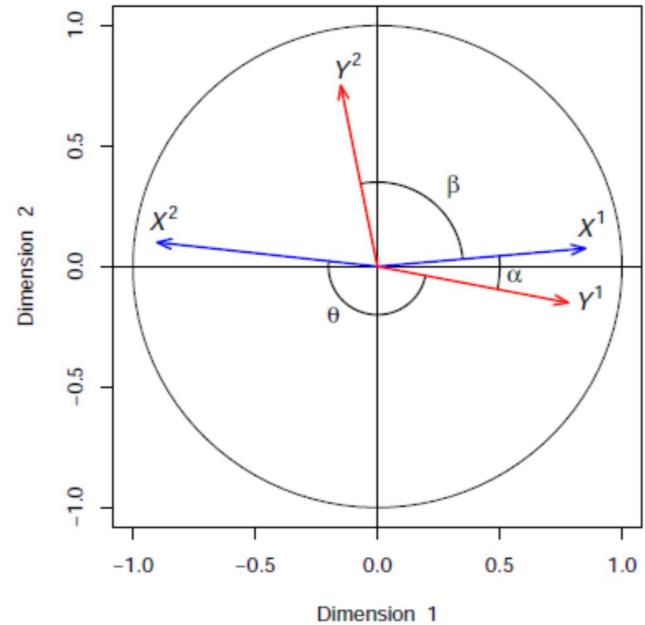
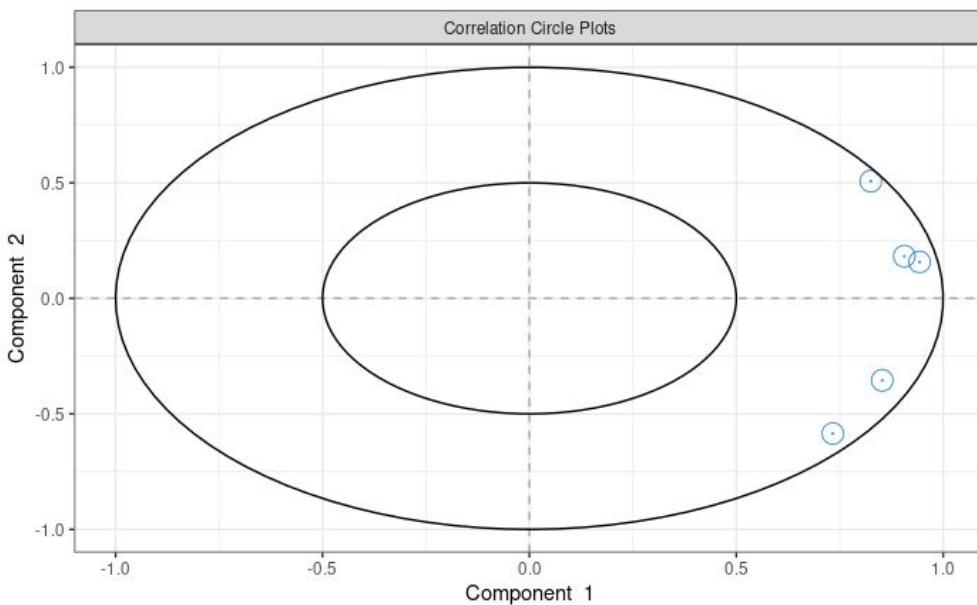


- 2 blocs
- 1 flèche = 1 échantillon.
- Origine = l'emplacement de l'échantillon en X, la pointe = l'emplacement dans Y
- Courtes flèches = X et Y sont en accord pour cet échantillon et inversement.

> `plotArrow(...)`

2. Variables plots

Circle Correlation plot

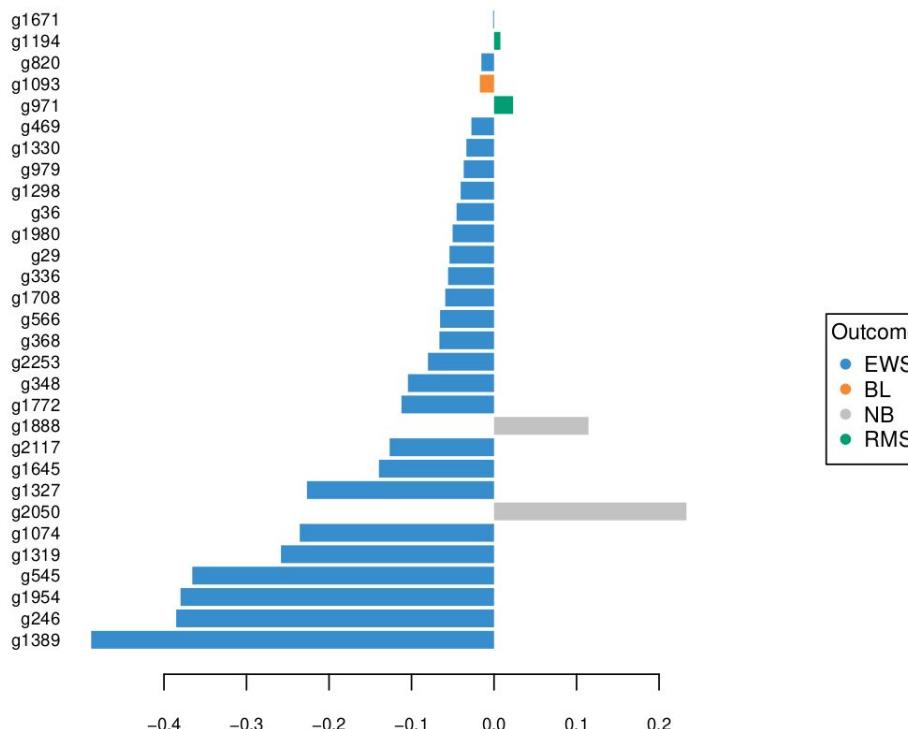


- Visualise les relations entre les variables et les composantes (*structure latentes*)
- Importance de la longueur du vecteur et du cosinus de l'angle avec l'axe
- Montre l'importance des variables / contribution par rapport aux composantes

> `plotVar(...)`

Loading plot

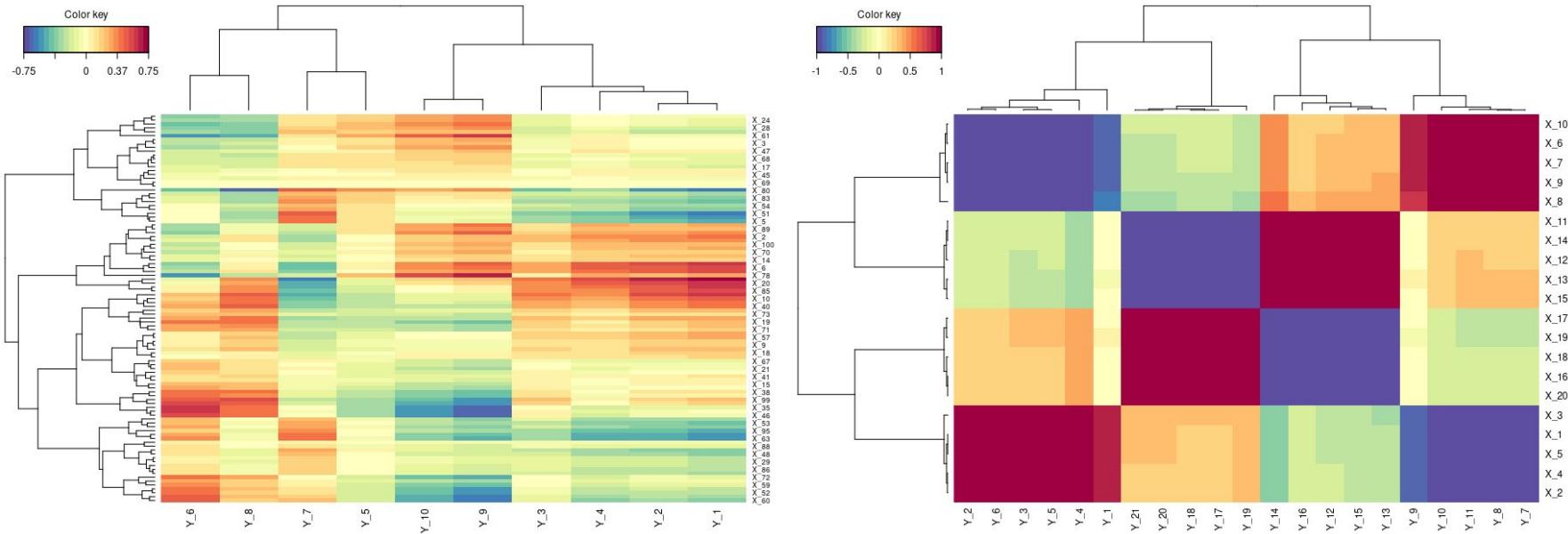
Contribution on comp 2



- Très similaire au CCplot
- Importance des variables chaque composantes (taille des barres)
- Ordonnée par ordre d'importance
- en mode supervisé, la couleur correspond au groupe le plus proche

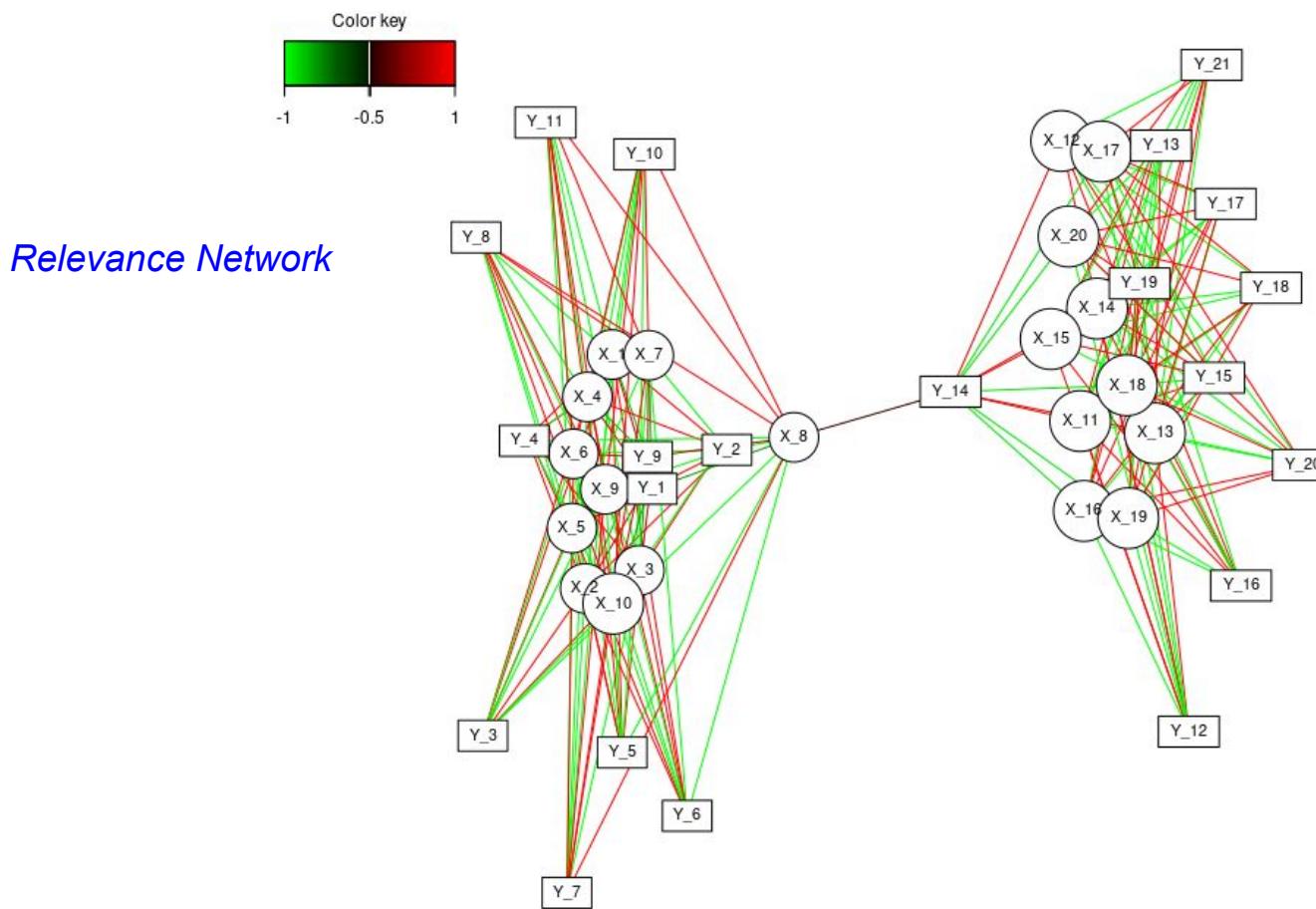
> `plotLoadings(...)`

Clustered Image Map



- Correlation = couleur
- Montre la corrélation entre des ensembles de variables de type différents (2 blocs)
- Clusters
- Dendrogramme

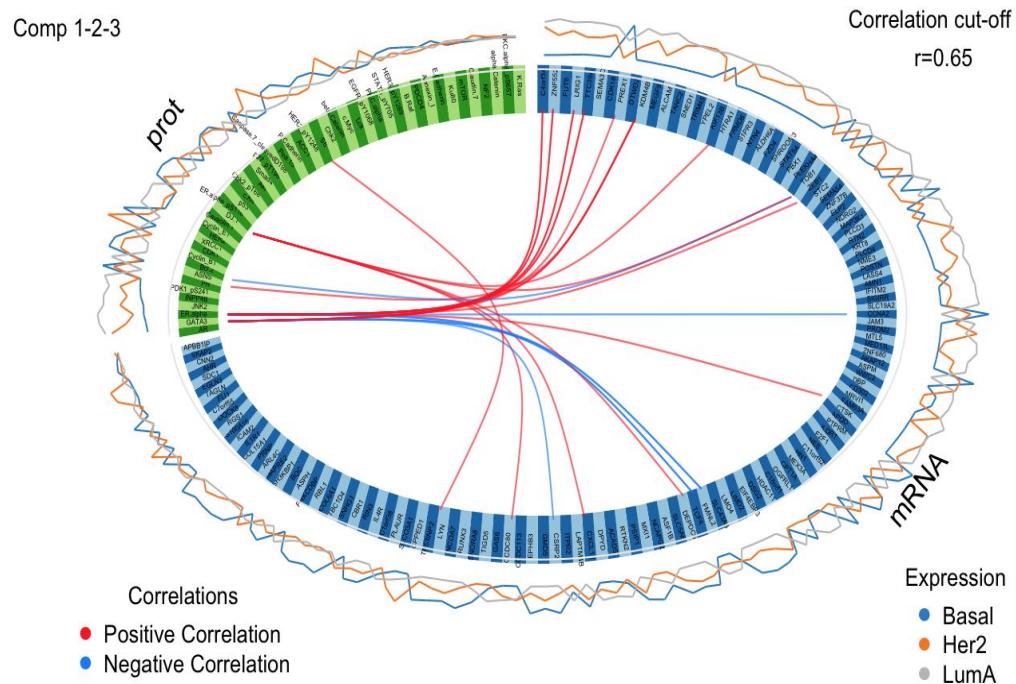
> `cim(...)`



- Montre la corrélation entre 2 variables (2 à 2) par un lien (tous les blocs).
- Corrélation = couleur du lien (*positif/négatif, fort/faible*)
- choix du **cutoff** (*affichage des liens*)
- Recherche de clique / sous-réseau

> `network(...)`

Circos plot



- Montre les associations / corrélation entre les blocs
- Corrélation = lien au centre (cutoff)
- Bordure externe = Intensité de l'expression (*par groupe en supervisé*)

> circos(...)

Mise en contexte

- omique ? multi-omique ?
- Pourquoi ? Difficultés
- Ressources

Méthodes d'intégration

- Intégration Conceptuelle
- Intégration Statistique
- Intégration à base de modèle

Méthodes Multivariées

- Analyse en Composante Principale
- Autres méthodes multivariées

Méthodes à base de réseaux biologiques

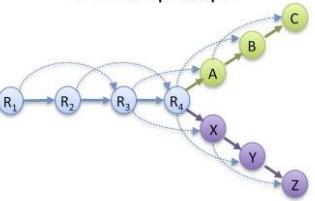
- Les réseaux en biologie
- Exploration des réseaux

Les Réseaux en Biologie

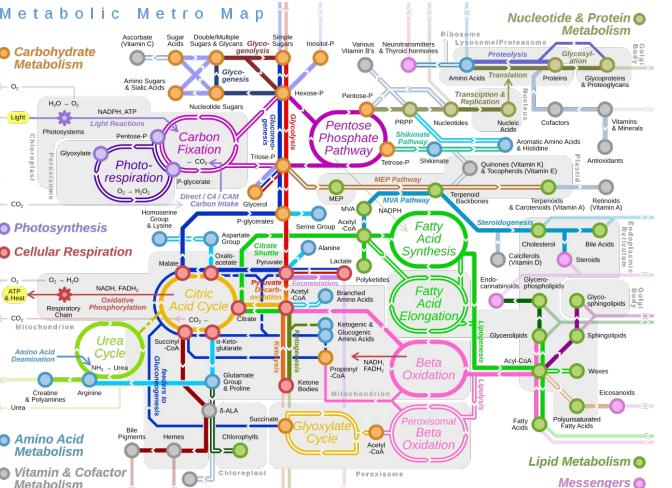
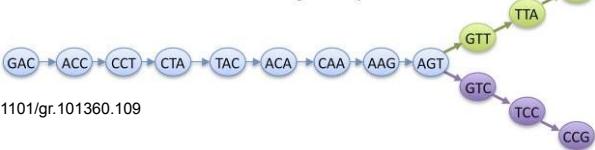
A Read Layout

```
R1: GACCTACA  
R2: ACCTACAA  
R3: CCTACAA  
R4: CTACAAGT  
A: TACAAAGTT  
B: ACAAGTTA  
C: CAAGTTAG  
X: TACAAGTC  
Y: ACAAGTCC  
Z: CAAGTCGG
```

B Overlap Graph



C de Bruijn Graph

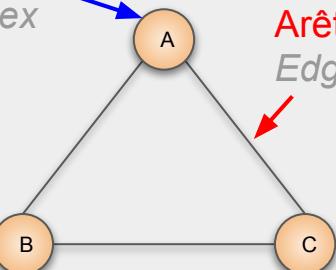


Les Graphes

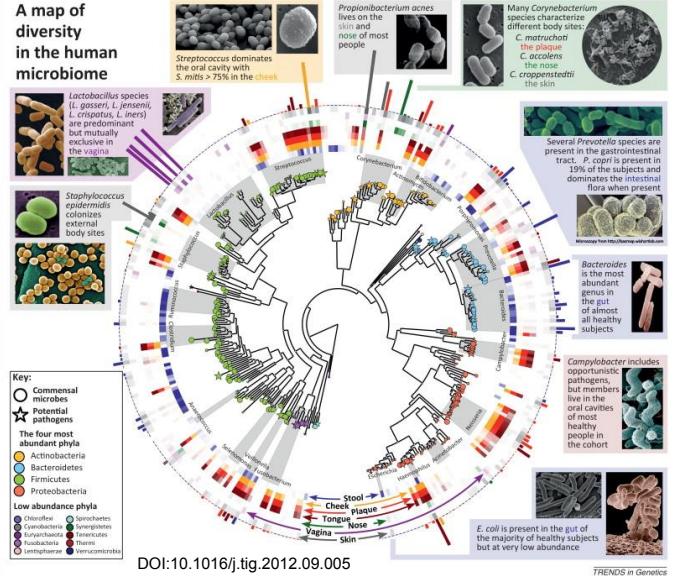
Noeud
Vertex

Arête
Edge

Graphe non-orienté:
Graphe orienté:



A map of diversity in the human microbiome



Gene Co-Expression / Regulation Network

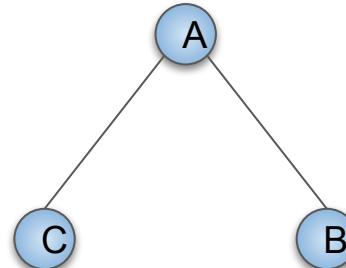
	Gene-A	Gene-B	Gene-C
Ech1	10	5	6
ECh2	42	49	2
...			



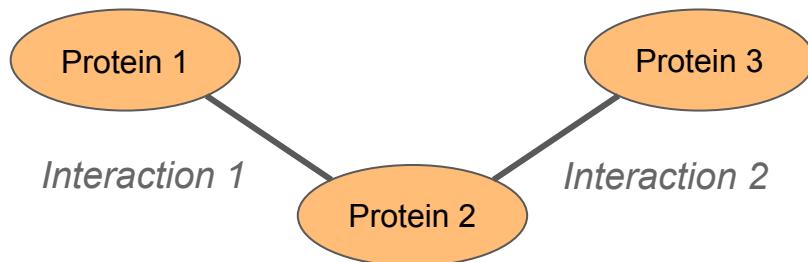
Matrice d'adjacence

	Gene-A	Gene-B	Gene-C
Gene-A	1	1	1
Gene-B	1	1	0
Gene-C	1	0	1

1: connecté 0: pas connecté



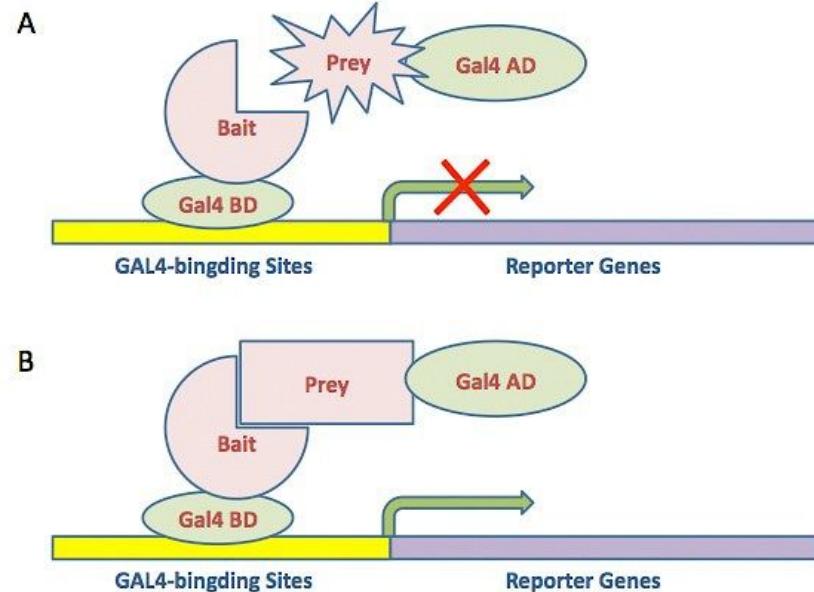
Réseaux d'interaction Protéine-Protéine (PPI)



Interaction = (activation, binding to, phosphorylation...)

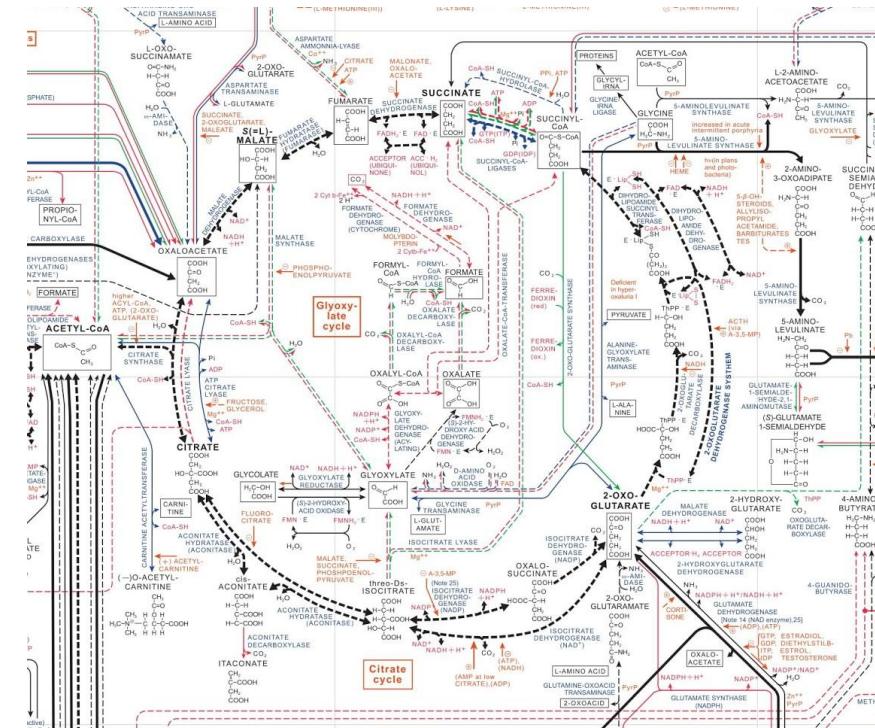
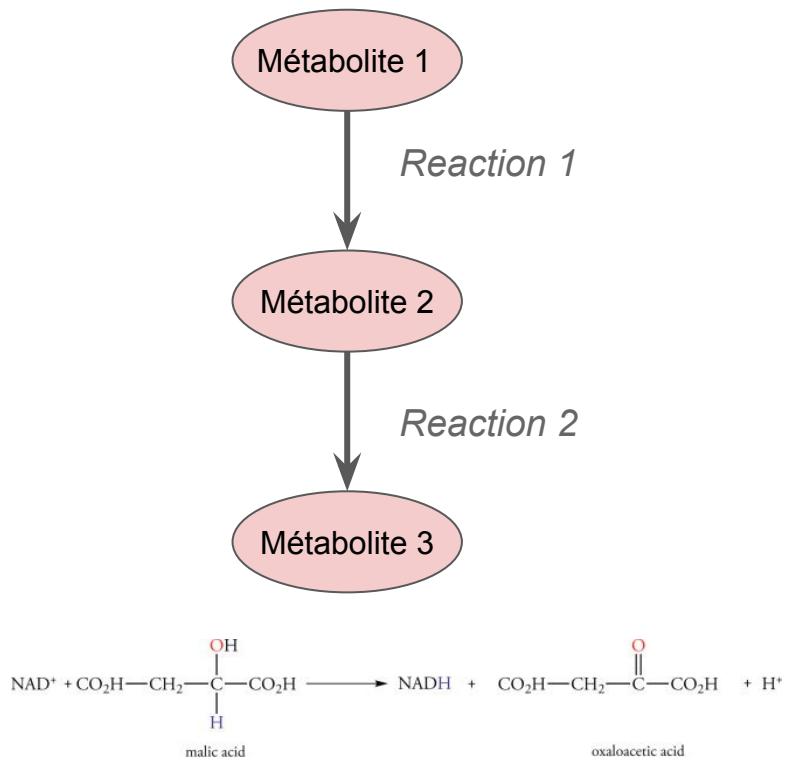


Two-hybrid screening



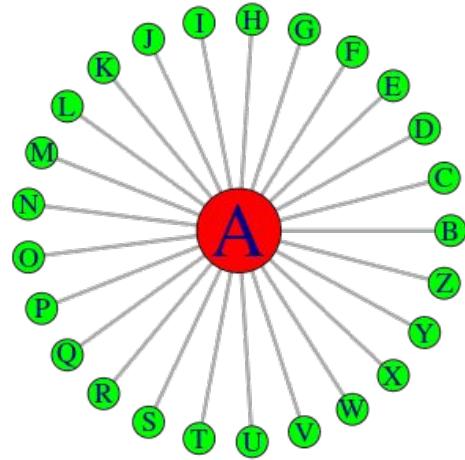
<http://www.myfolio.com/art/fez57ey2az>

Metabolic Pathways



Interaction = (synthesis, oxydation, hydration, decarboxylation, ...)

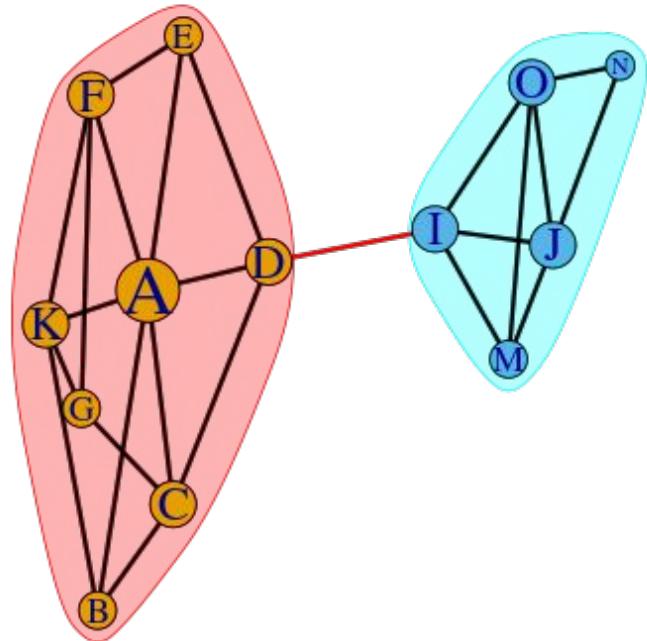
Analyse de Topologie



❑ Degree distribution

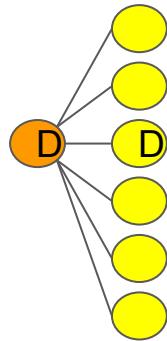
- ❑ Modularity
- ❑ Hierarchy
- ❑ Centrality
- ❑ ...

❑ Shortest path
❑ Clustering



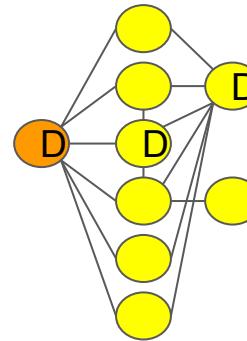
Analyse de Topologie

Random Walk



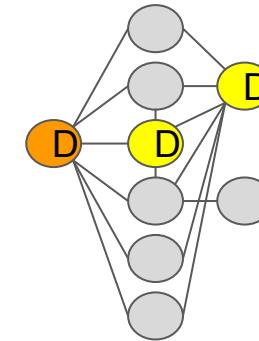
Direct neighbors

- False positives
 - False negatives



Shortest Path

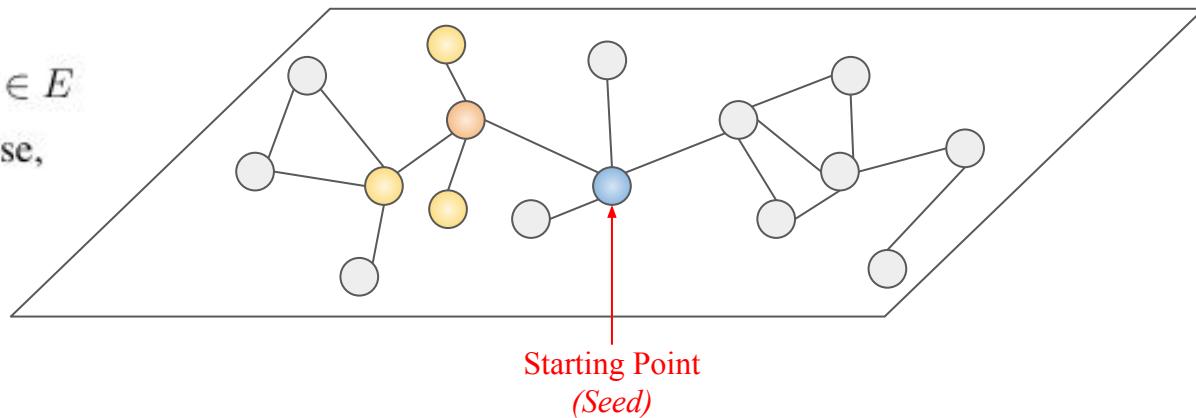
- *Small world*
 - False positives



Network Propagation

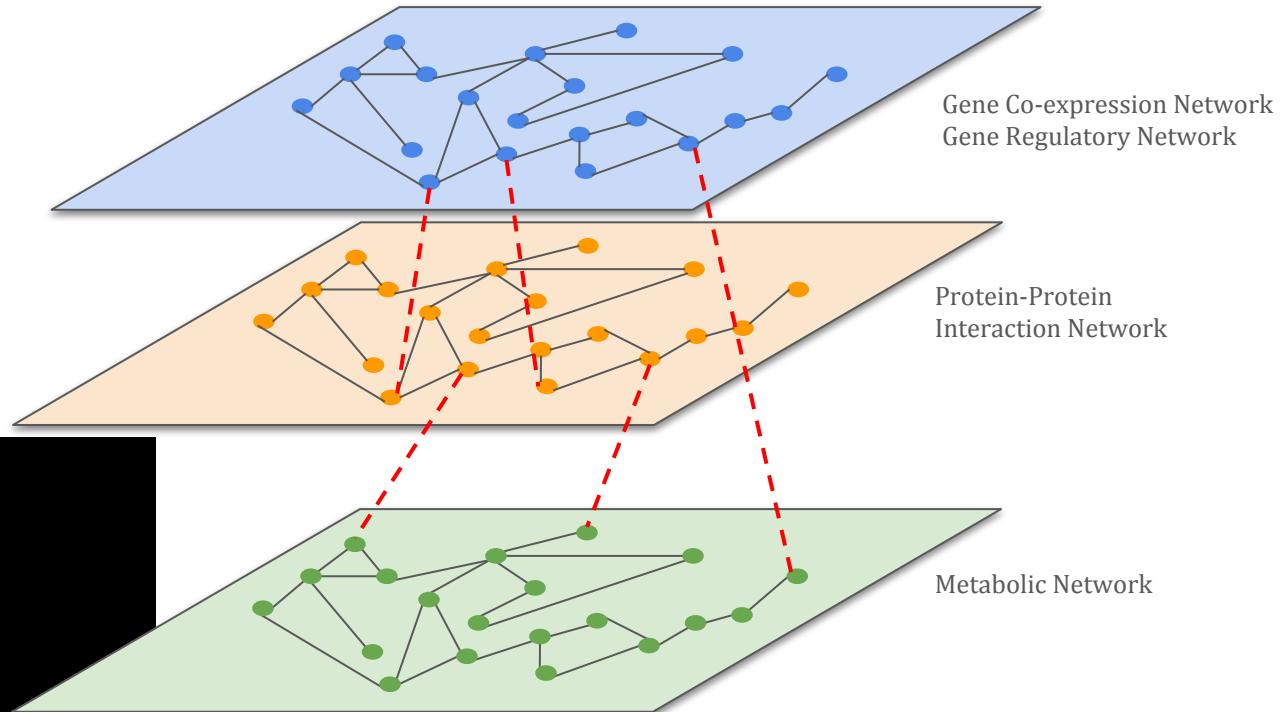
- Random Walk
 - Node Ranking

$$\mathbb{P}(v_{t+1} = y | v_t = x) = \begin{cases} \frac{1}{d(x)} & \text{if } (x, y) \in E \\ 0 & \text{otherwise,} \end{cases}$$



Multi-Omics Networks

Multi-Layered



Multiplex Heterogeneous Network

Mise en contexte

- omique ? multi-omique ?
- Pourquoi ? Difficultés
- Ressources

Méthodes d'intégration

- Intégration Conceptuelle
- Intégration Statistique
- Intégration à base de modèle

Méthodes Multivariées

- Analyse en Composante Principale
- Autres méthodes multivariées

Méthodes à base de réseaux biologiques

- Les réseaux en biologie
- Exploration des réseaux

Dernières tendances

❑ Omiques inhabituels

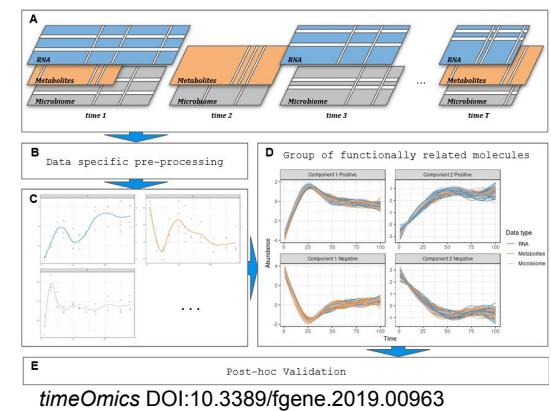
Microbiome: iHMP (<https://www.hmpdacc.org/ihmp/>)



❑ Aspect Longitudinale

Suivre la dynamique des systèmes (complexe) au niveau moléculaire

❑ Méthode d'exploration des graphes (Random Walk, ...)

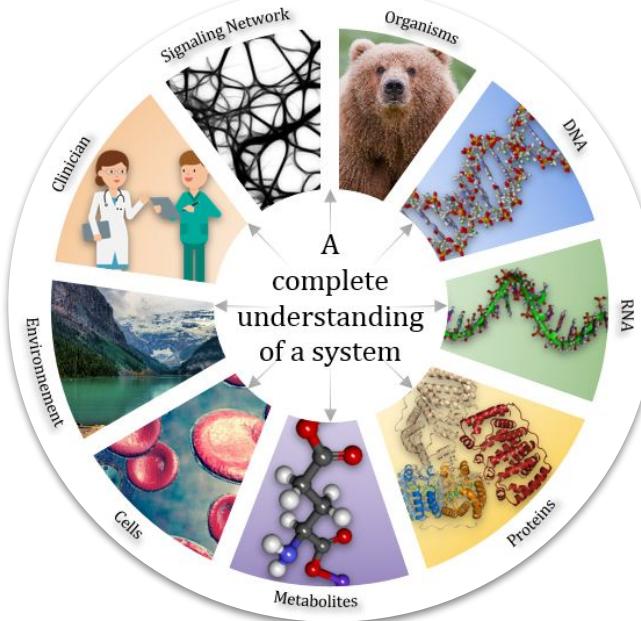
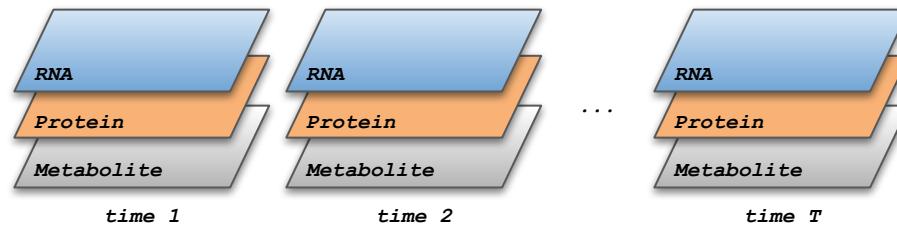


Intégration longitudinale



New experimental **design**

Multi-OMICS **time** series data



Intégration longitudinale

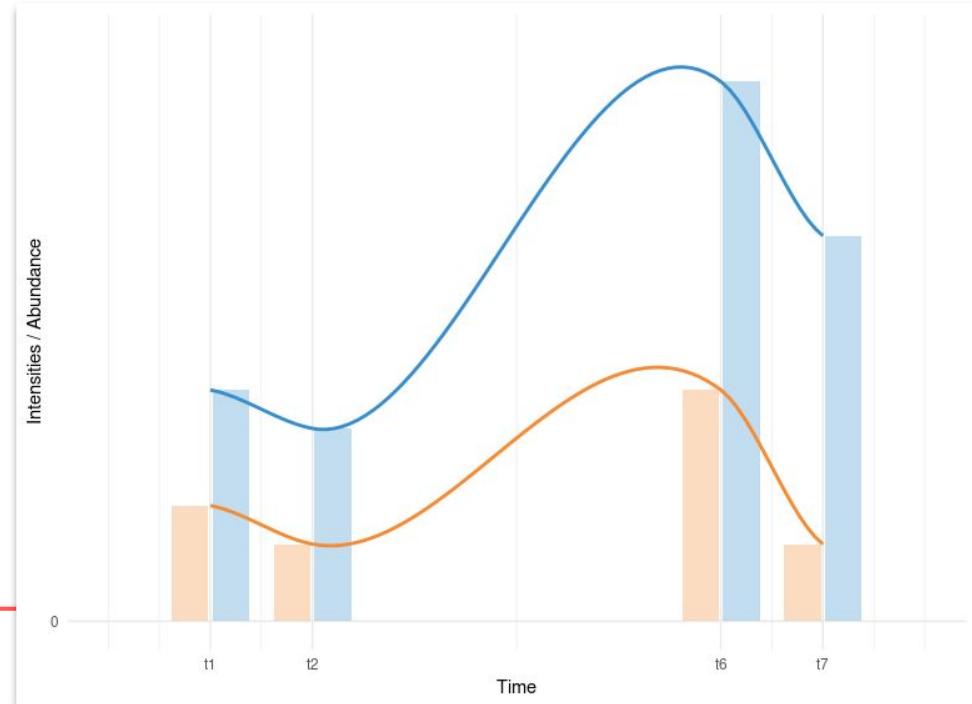
- Time effect
 - ◆ between two time points
 - ◆ between intervals of time points

- Group effect

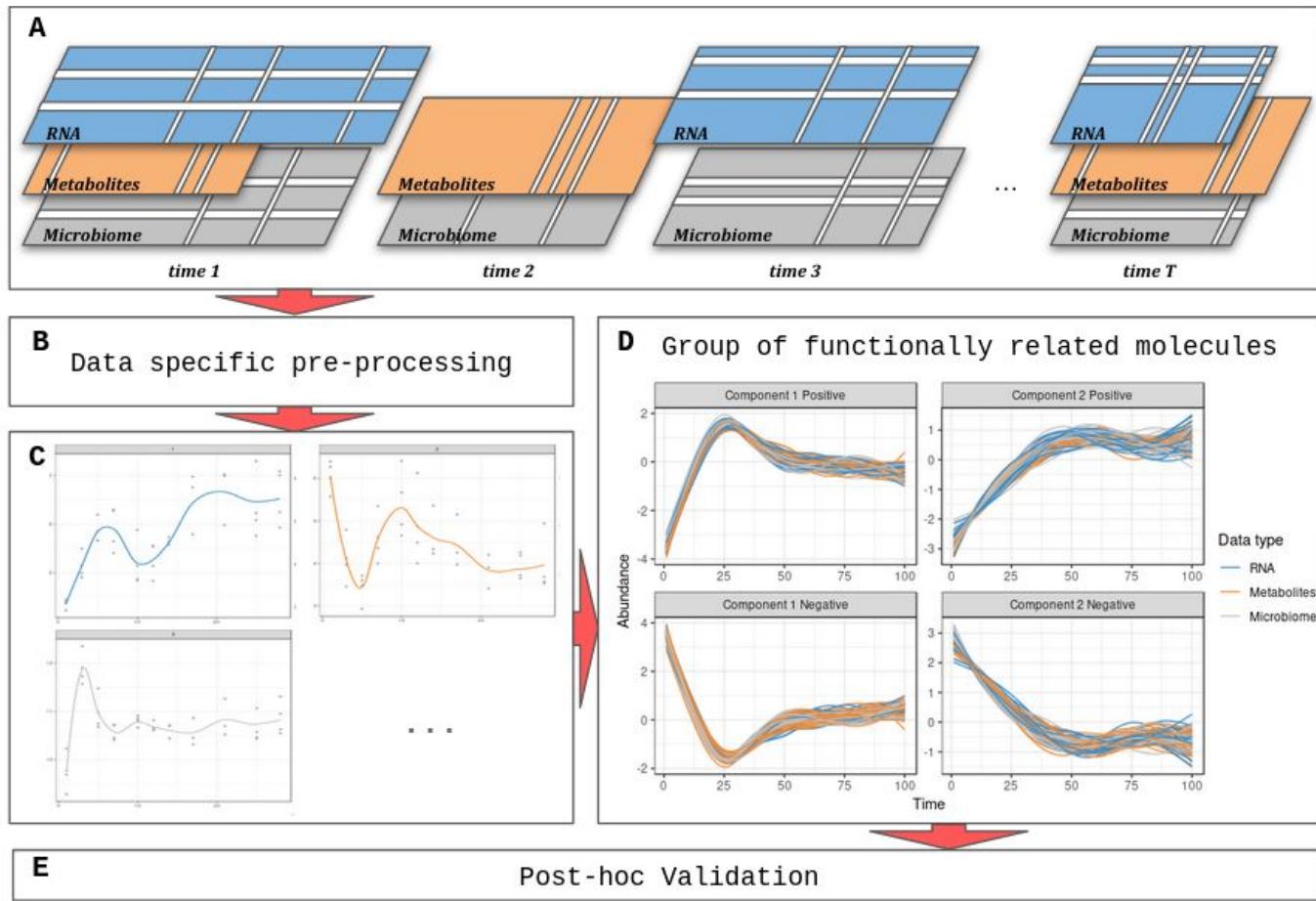
- Detect similar pattern

Why cluster data ?

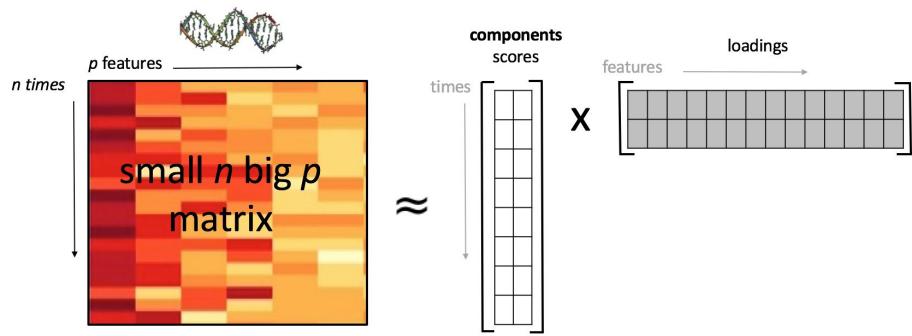
- Group molecules with similar expression profile
- Functionally related
- Can be used to infer regulatory relationships



- A.** Data Acquisition
- B.** Pre-processing
- C.** Modelling and Filtering
- D.** Clustering / Integration
- E.** Validation



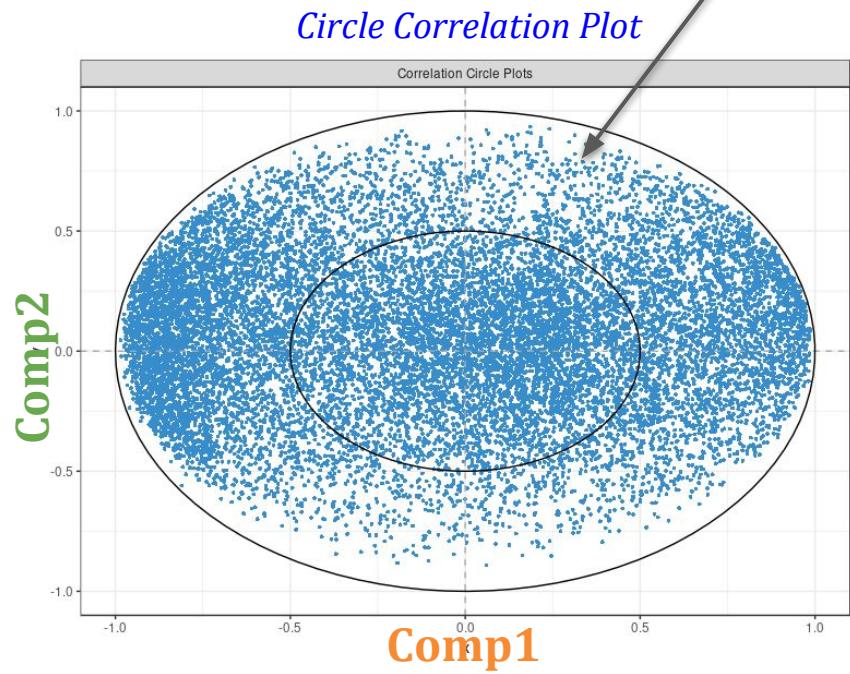
Principal Component Analysis - PCA



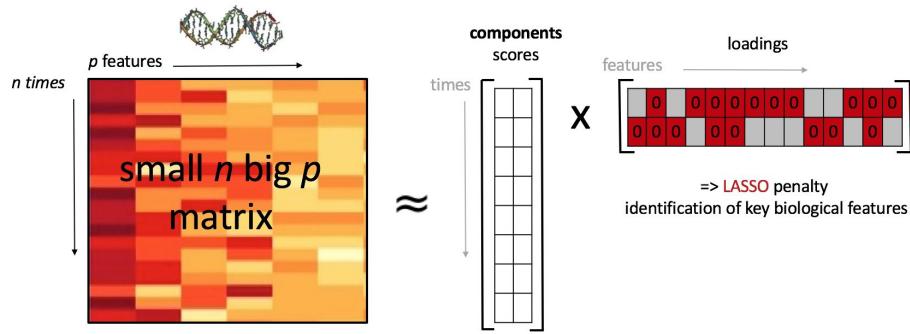
- Block = OTU table
- Unsupervised
- Dimension reduction
- Maximize **variance**

([mixOmics](#) R package)

1 point = 1 time profile (genes, proteins, ...)



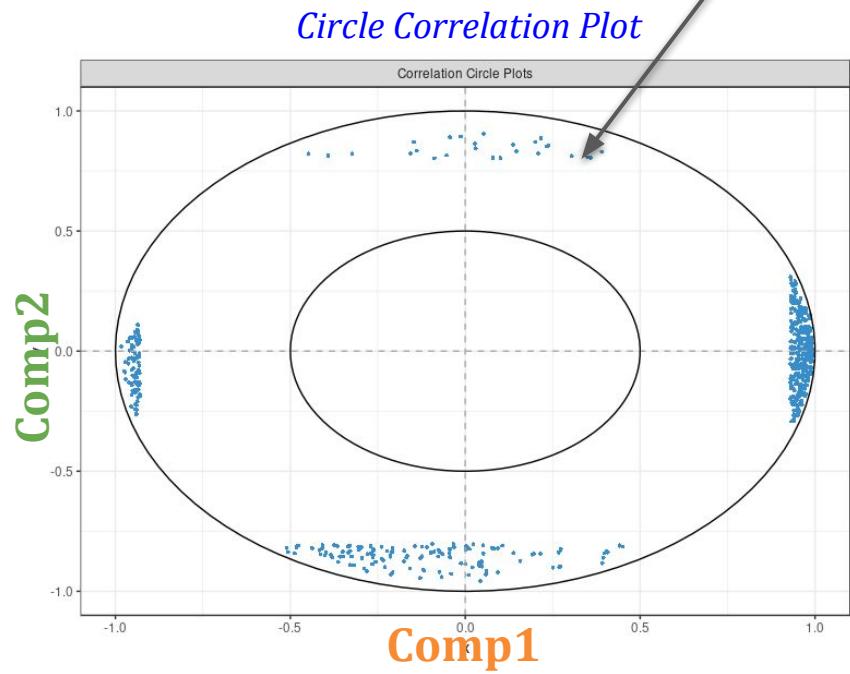
sparse Principal Component Analysis - sPCA



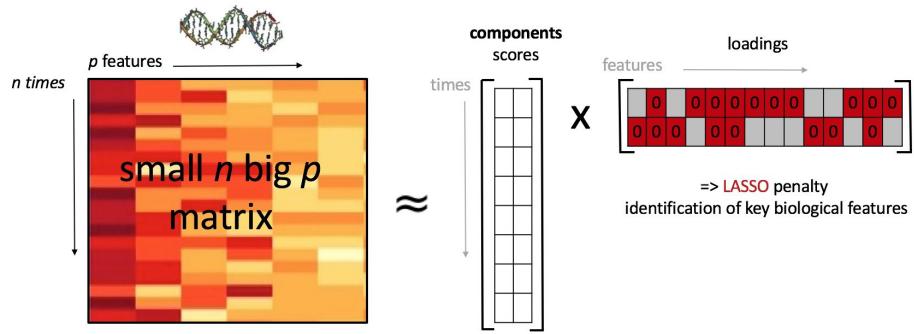
- Block = OTU table
- Unsupervised
- Dimension reduction
- Maximize **variance**
- Sparse version for **feature selection**

([mixOmics](#) R package)

1 point = 1 time profile (genes, proteins, ...)



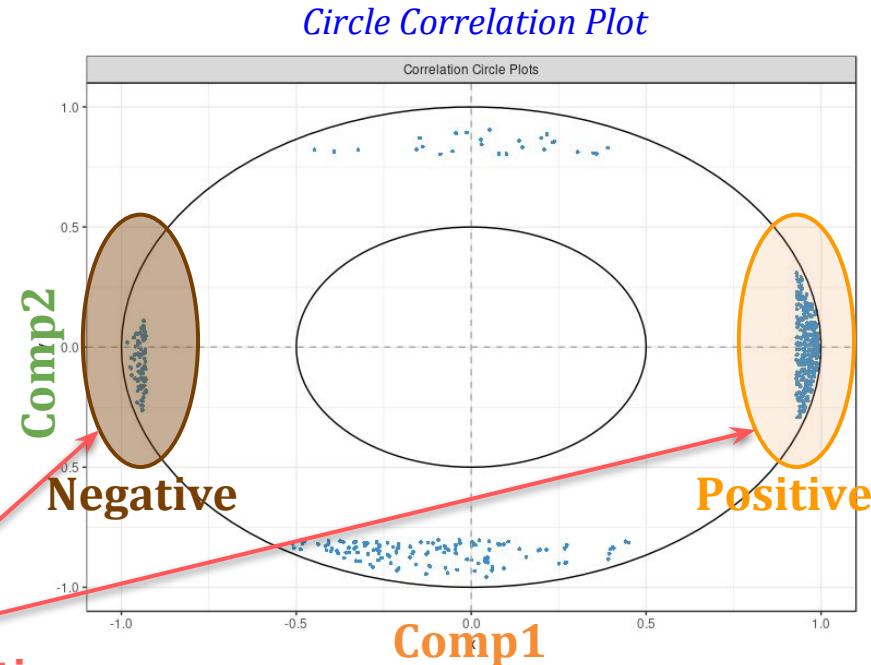
sparse Principal Component Analysis - sPCA



Each component selects **2 clusters of OTUs**

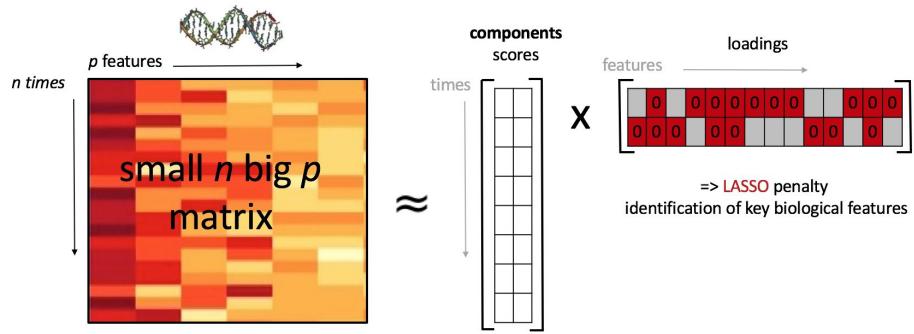
- Positively correlated
- Negatively correlated

Negative Correlation



Two clusters identified on **comp 1**

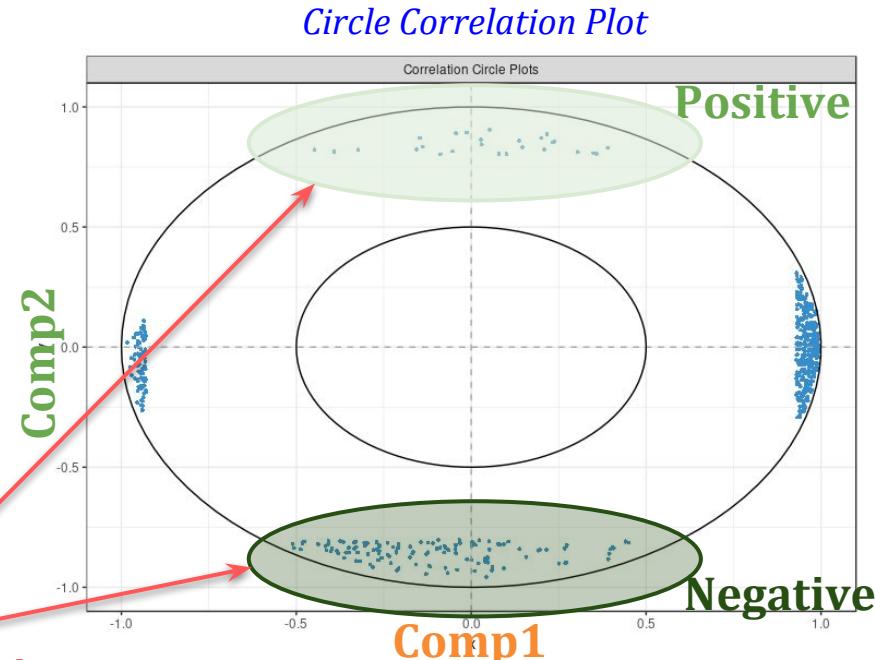
sparse Principal Component Analysis - sPCA



Each component selects **2 clusters of OTUs**

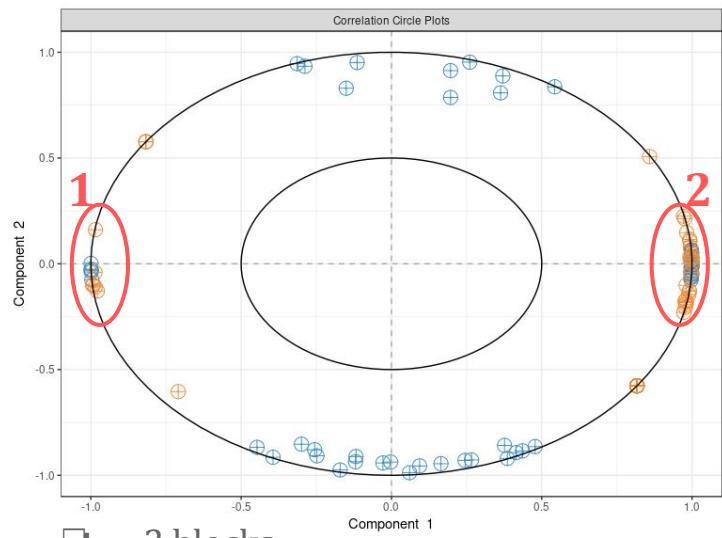
- Positively correlated
- Negatively correlated

Negative Correlation

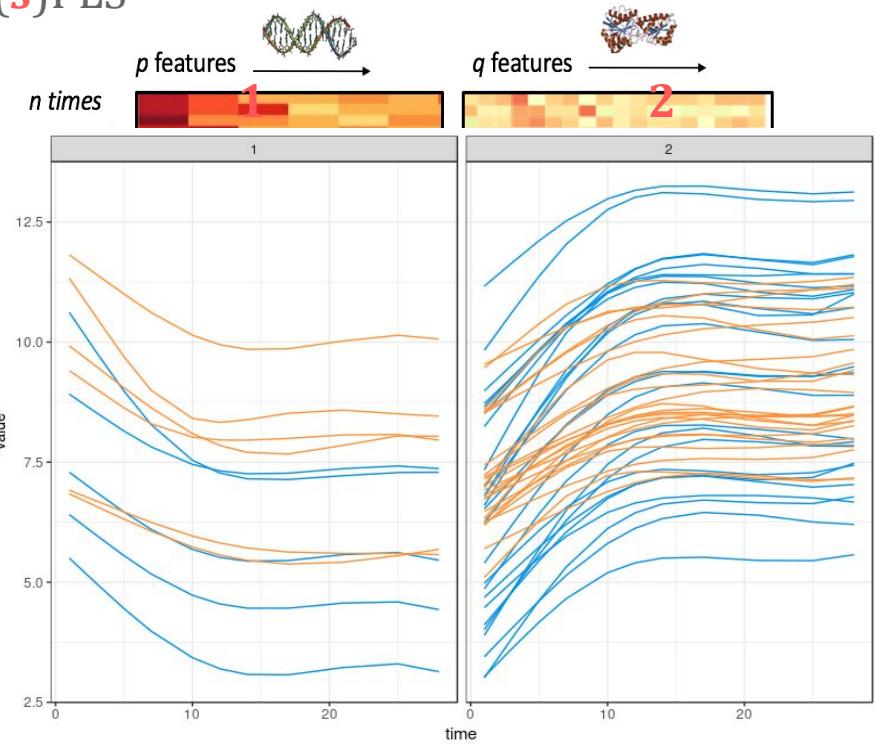


Two clusters identified on **comp 2**

(sparse) Projection on Latent Structures - (s)PLS



(mixOmics R package)

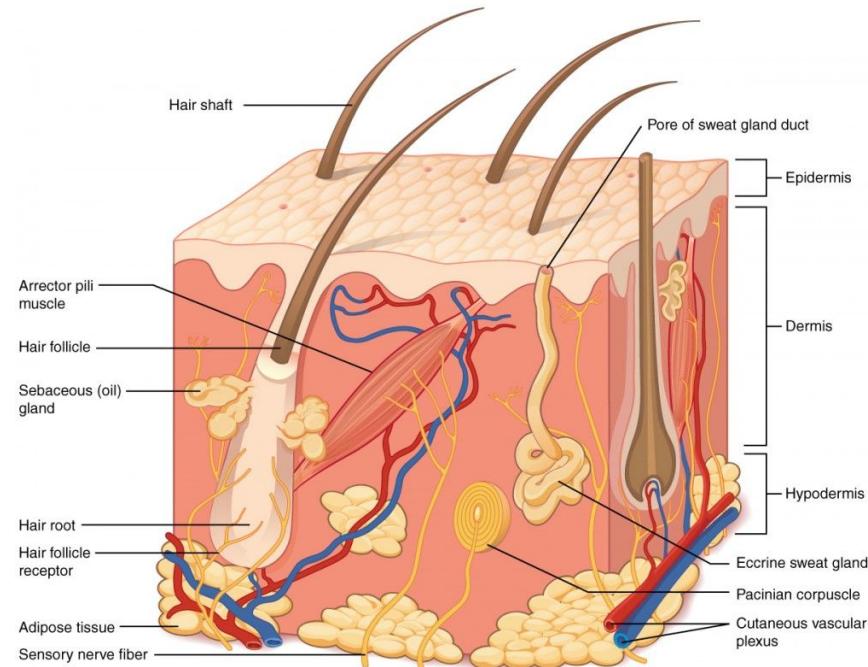


Example: *epidermis reconstruction*

Main function of the human skin:

- ❑ Organ of Sensation
- ❑ Thermo-regulation
- ❑ Secretion
- ❑ Immunité
- ❑ Psychological
- ❑ **Protection**

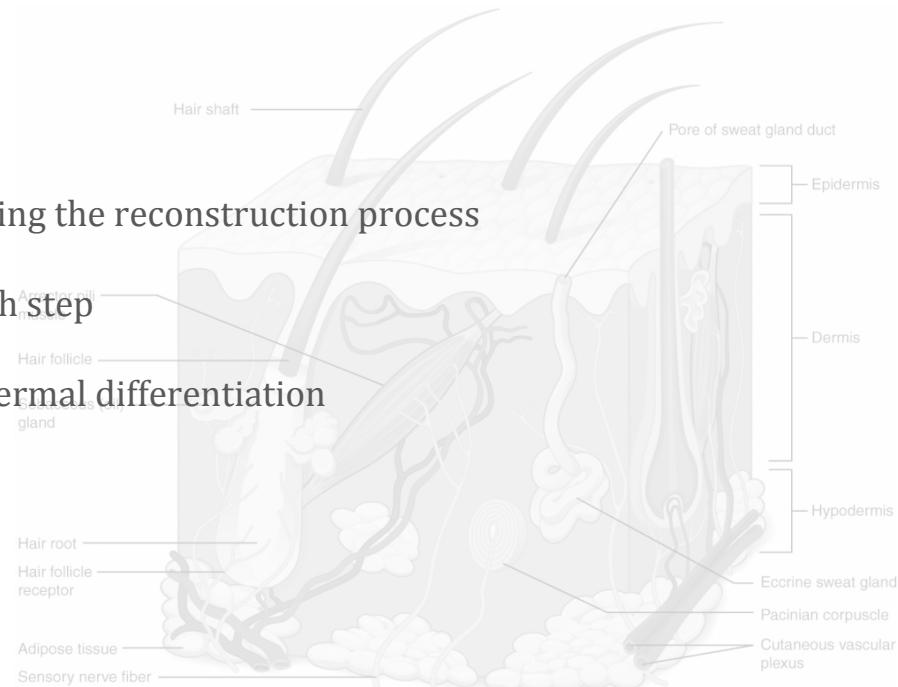
High capacity for regeneration



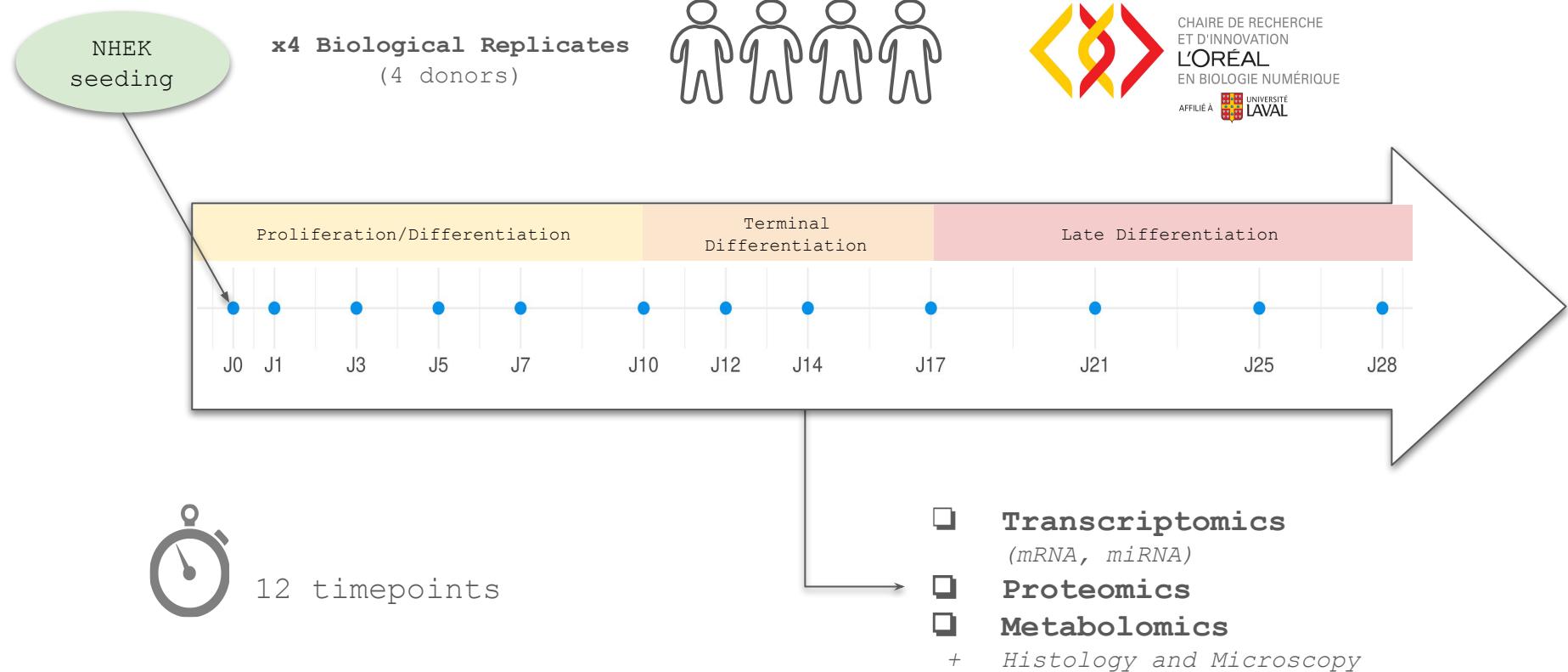
We aim to study the kinetics of **epidermal differentiation** during the **reconstruction process**.

Example: *epidermis reconstruction*

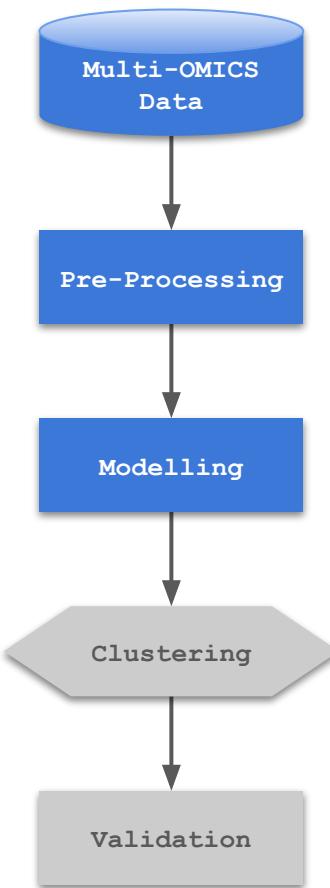
- Understand the epidermal differentiation during the reconstruction process
- Identify key biological players involved in each step
- Identify key and transitional steps of the epidermal differentiation
- Identify key missing players



Example: *epidermis reconstruction*



Example: *epidermis reconstruction*



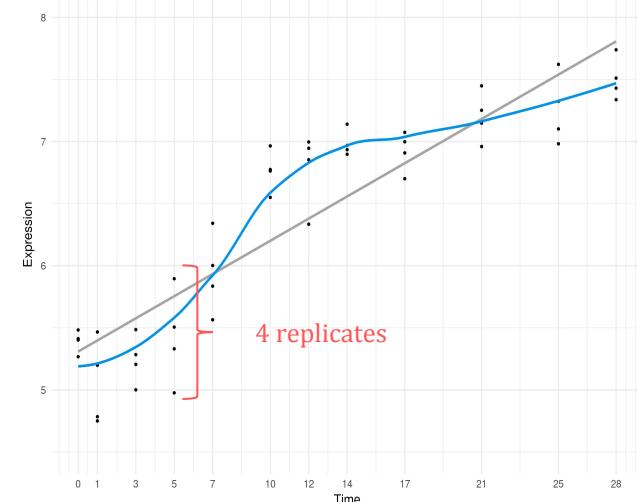
OMICS	# Molecules	# Molecule after filtering
mRNA	15,348	1,979 (13%)
miRNA	896	171 (19%)
Proteomics	1,820	987 (54%)
Metabolomics	428	338 (79%)

Filtering

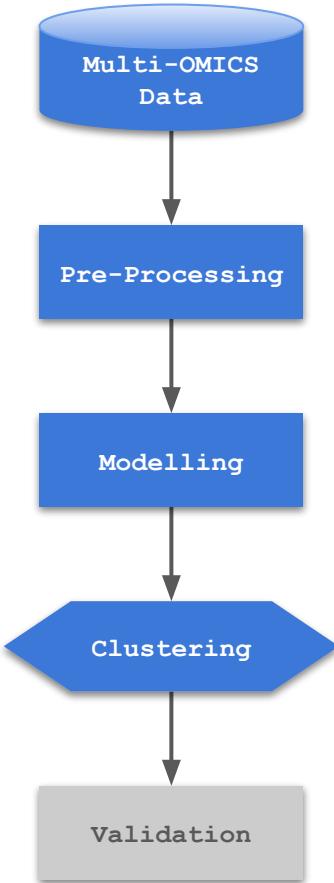
- “flat” profiles
- Molecules with no D.E. in time

Modelling

- Linear Mixed Model Spline
- Filter noisy profile



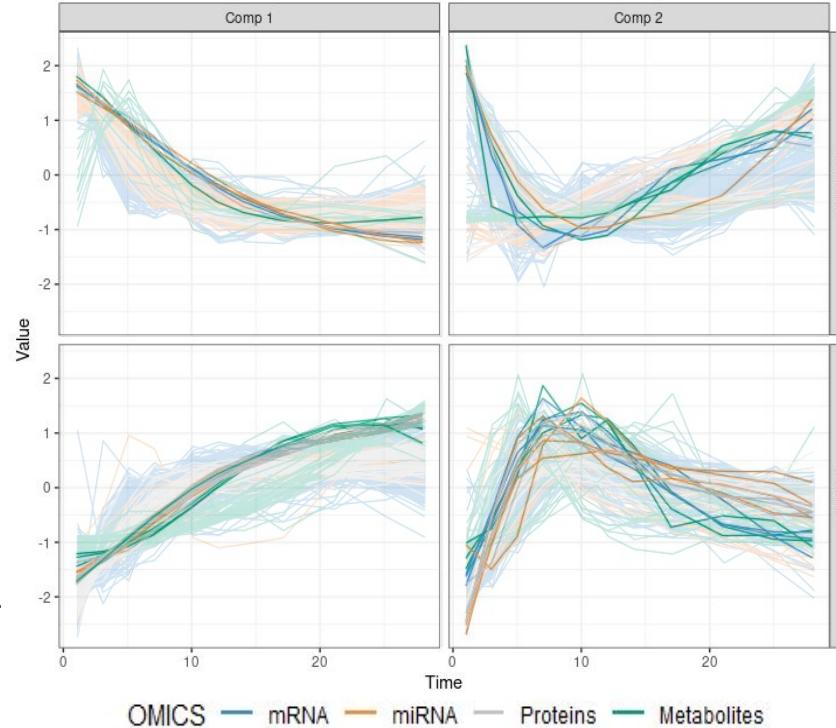
Example: *epidermis reconstruction*



Cluster -1 (*Comp1 Neg*)

mRNA: 418 (5)
miRNA: 87 (18)
Proteins: 179 (5)
Metabolites: 62 (16)

Multi-OMICS (sparse) clusters



Cluster 1 (*Comp1 Pos*)

mRNA: 898 (10)
miRNA: 21 (17)
Proteins: 317 (20)
Metabolites: 127 (9)

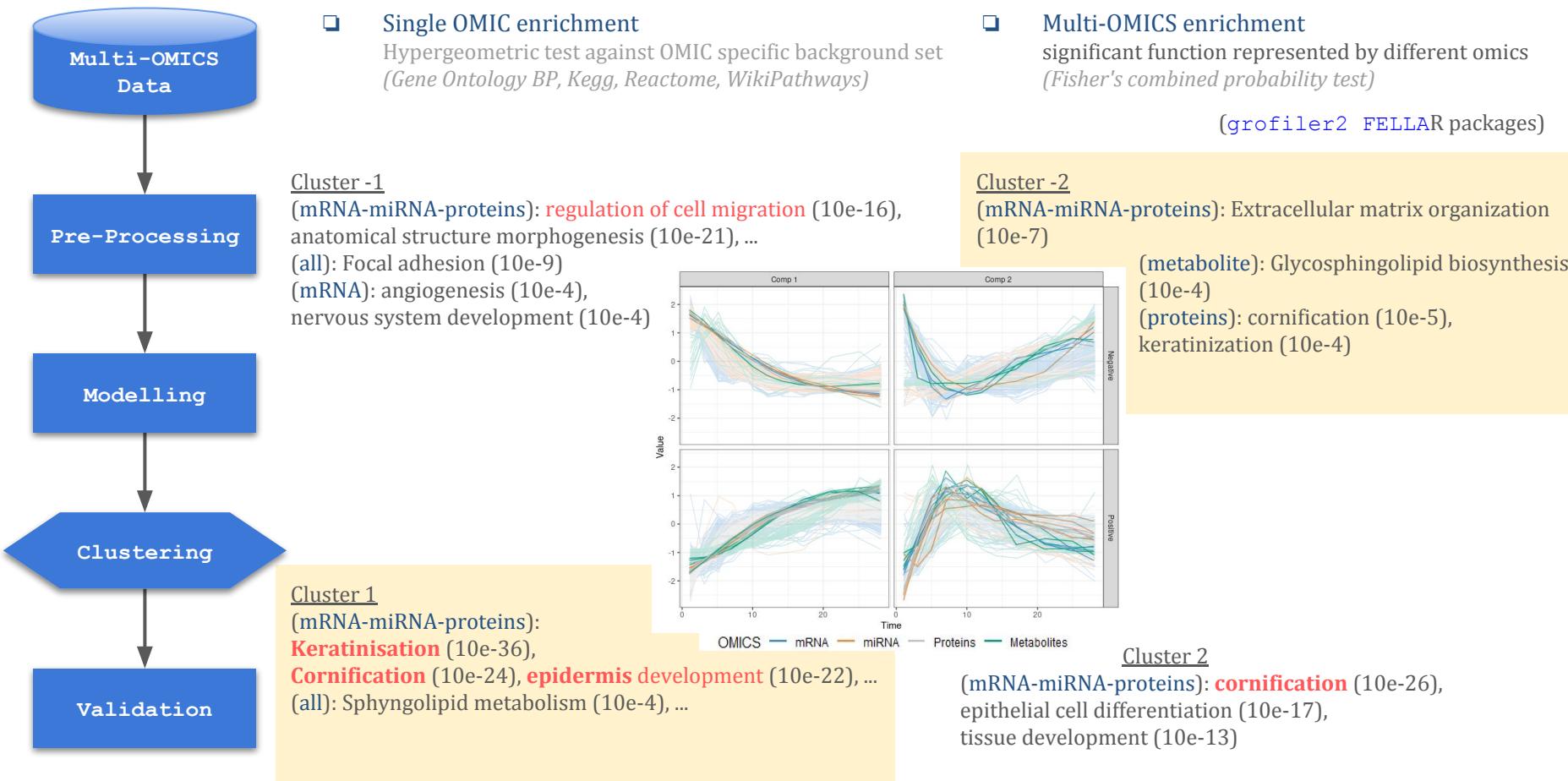
Cluster -2 (*Comp2 Neg*)

mRNA: 228 (7)
miRNA: 10 (7)
Proteins: 197 (3)
Metabolites: 58 (6)

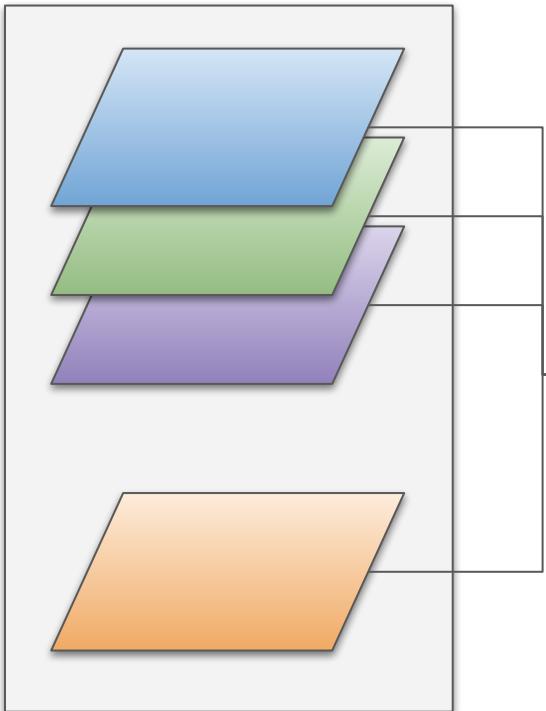
Cluster 2 (*Comp2 Pos*)

mRNA: 435 (8)
miRNA: 53 (8)
Proteins: 294 (12)
Metabolites: 91 (9)

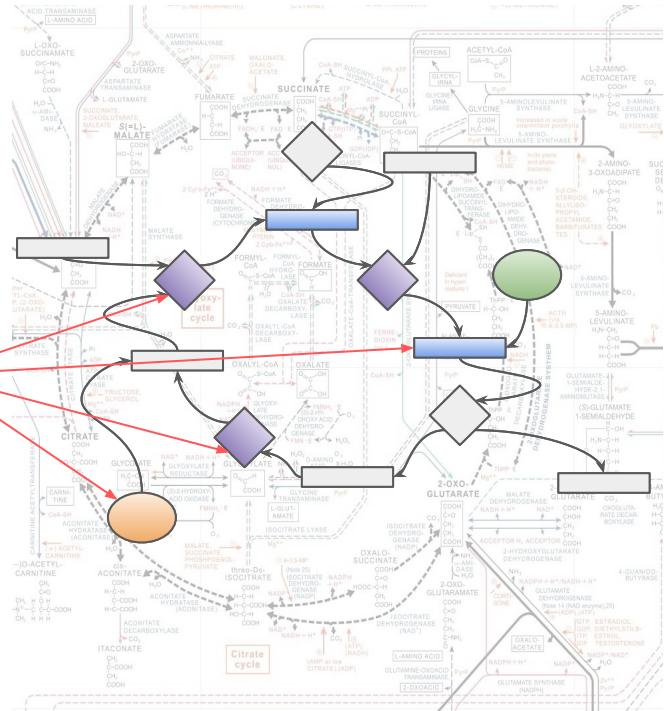
Example: *epidermis reconstruction*



Example: *epidermis reconstruction*

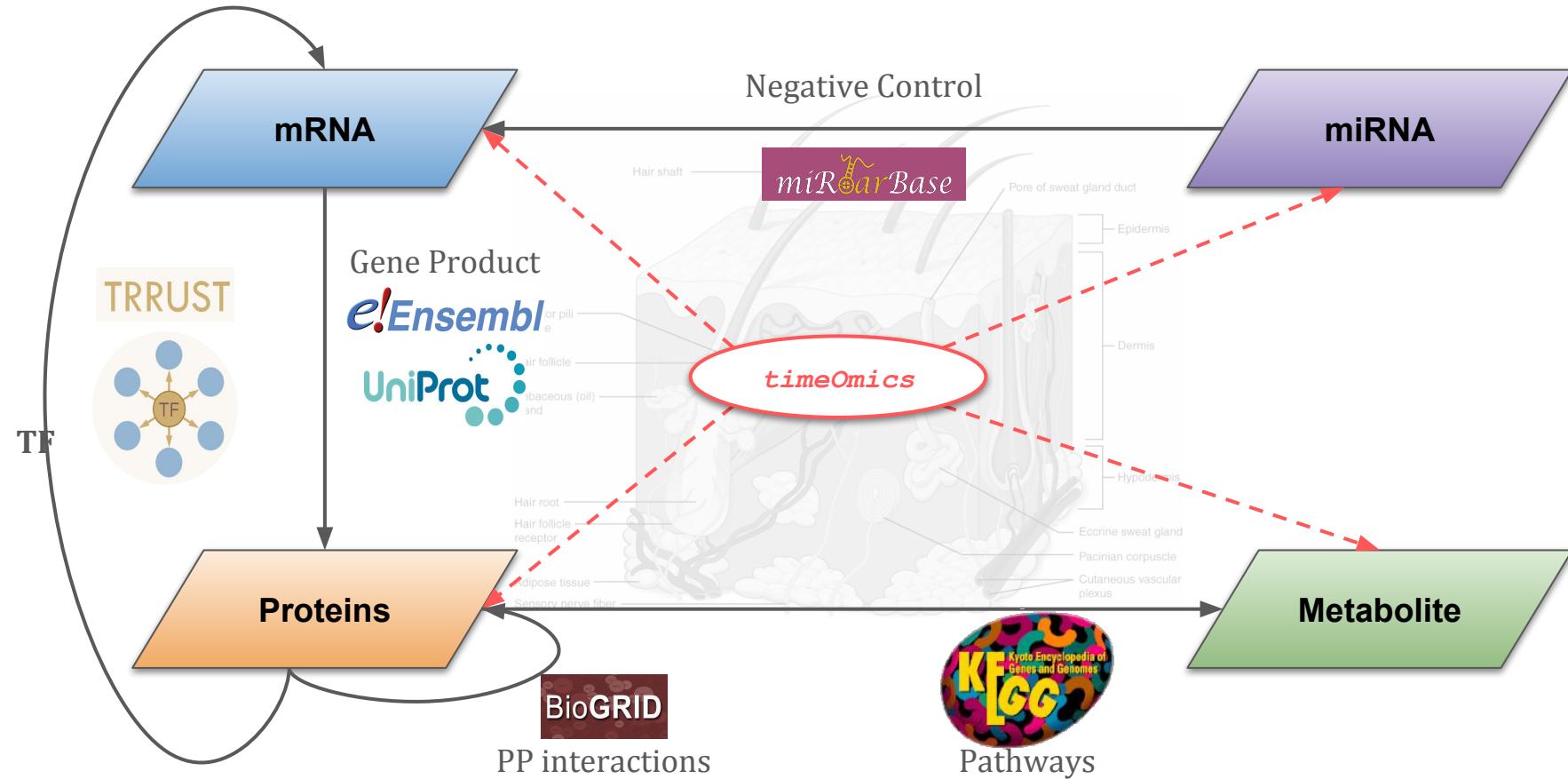


Multi-OMICS
Signature



Knowledge

Example: *epidermis reconstruction*



- <http://melgen.org/multi-omics-approach/>
- Huang S, Chaudhary K, Garmire LX. More Is Better: Recent Progress in Multi-Omics Data Integration Methods. *Frontiers in Genetics*. 2017;8:84. doi:10.3389/fgene.2017.00084.
- Hasin, Y., Seldin, M., & Lusis, A. (2017). Multi-omics approaches to disease. *Genome biology*, 18(1), 83.
- Cavill, R., Jennen, D., Kleinjans, J., & Briedé, J. J. (2015). Transcriptomic and metabolomic data integration. *Briefings in bioinformatics*, 17(5), 891-901.
- Wanichthanarak, K., Fahrmann, J. F., & Grapov, D. (2015). Genomic, proteomic, and metabolomic data integration strategies. *Biomarker insights*, 10(Suppl 4), 1.
- <https://www.frontiersin.org/research-topics/2280/multi-omic-data-integration>
- <http://mixomics.org>
- L'analyse en Composante Principale, Sébastien Déjean - math.univ-toulouse.fr