

# Computational biology laboratory CHU de Québec – Université Laval

## Current Projects

Arnaud Droit, PhD



CHAIRE DE RECHERCHE  
ET D'INNOVATION  
L'ORÉAL  
EN BIOLOGIE NUMÉRIQUE  
AFFILIÉ À UNIVERSITÉ  
LAVAL



UNIVERSITÉ  
**LAVAL**

# Québec



États-Unis

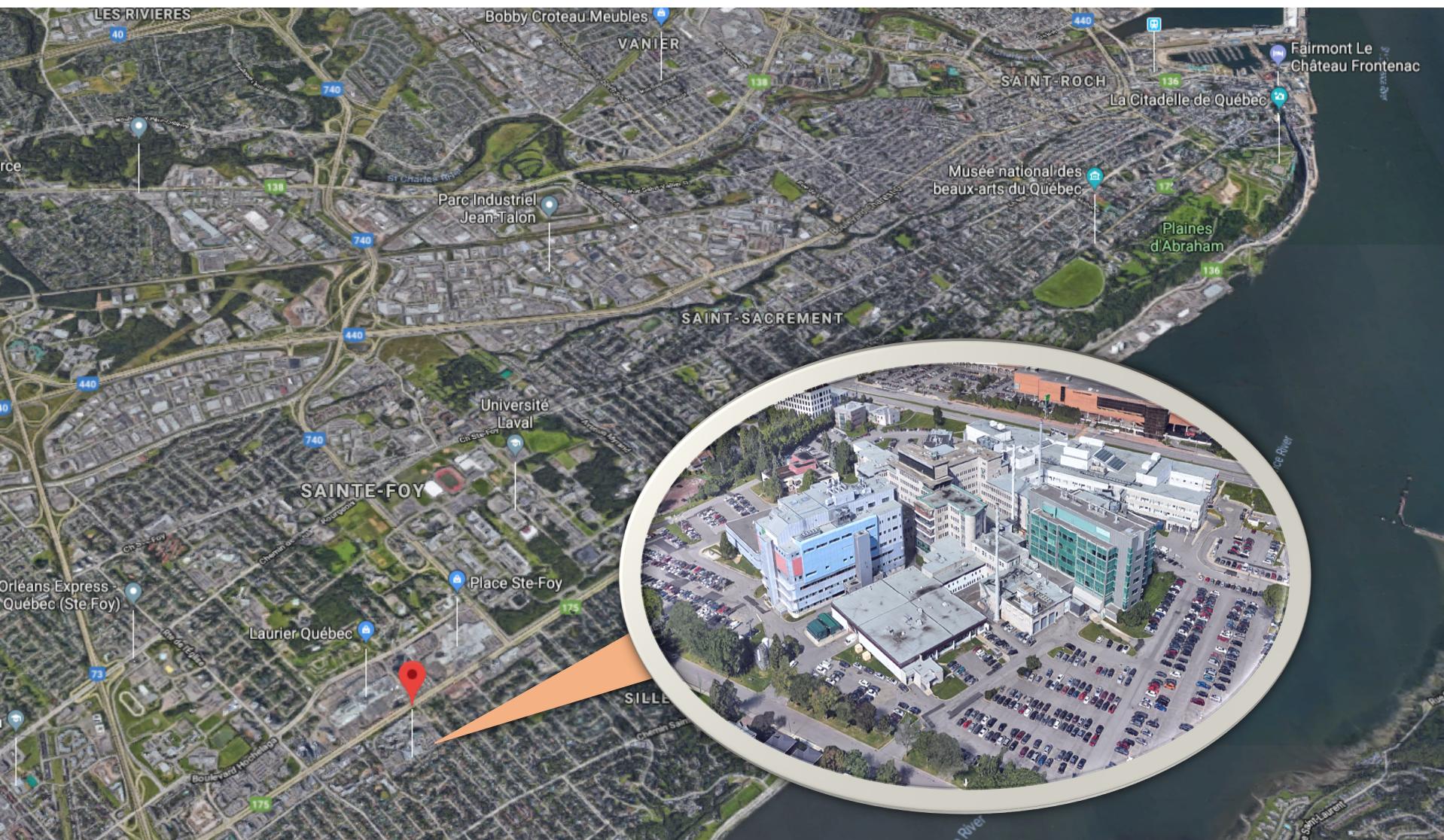
# Ville de Québec



# Université Laval – Ville de Québec



# Research Center - CHUL



# The Team (Not complete)



Mickael Leclercq, PhD



Christophe  
Tav  
(Doctorant)



Charles Joly Beauparlant, PhD



Eric  
Fournier



Julien Prunier,  
PhD



Kodjovi Dodji  
Mлага, PhD



Régis Ongaro  
(Doctorant)



Alban Mathieu



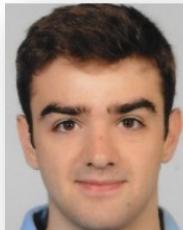
Loic Mangnier



Marie-Pier Scott Boyer,  
PhD



Frédéric  
Fournier



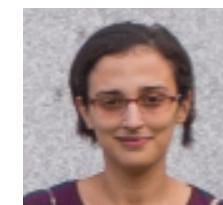
Antoine  
Bodein  
(Doctorant)



Francois Belleau



Afshin  
Jamshidi, PhD



Kawhla Seddiki  
(Doctorante)



Gwenaëlle  
Lemoine  
(Doctorante)

# Computational biology laboratory



- 8 research associates (Developers, Bioinformaticians, Biostatisticians, Data scientists...)
- 7 PhD students
- 2 Post-Doc
- some Interns
- Main projects: Multi-OMICS Integrations, Exome sequencing, Epigenetics, Proteomics, Automatic classification
- Linked to Compute Canada infrastructure



# Main activities and expertise

- We develop tools and strategies dedicated to
  - Multi-omics data analysis
  - Discovery of biomarker signatures using machine learning approaches
  - Elaboration of deep learning predictive models
  - Integration, exploration and visualization of biological big data
- We offer services to biologists through collaborations to analyse various types of omic data
  - Genomic (WGS, WES), Transcriptomic, Proteomic, Metabolomic, Epigenomic, Meta-genomic (microbiota)
- We have an expertise in the analysis of data in various diseases
  - Alzheimer, Multiple sclerosis, Hormonodependant cancers (Breast, Prostate), Leukemia, Fragile X syndrome, etc.
- L'Oréal Research and Innovation Chair in Digital Biology
  - Characterization of artificial skin reconstruction using Time-dependant multi-omics, including skin microbiota
- Currently in the Bioinformatics laboratory
  - 23 people (1 Professor, 7 Research associates, 6 Phds, 1 master, 2 postdoc, 5 interns, 1 sysadmin)
  - About 20 research projects

# Our infrastructure

## Computing resources

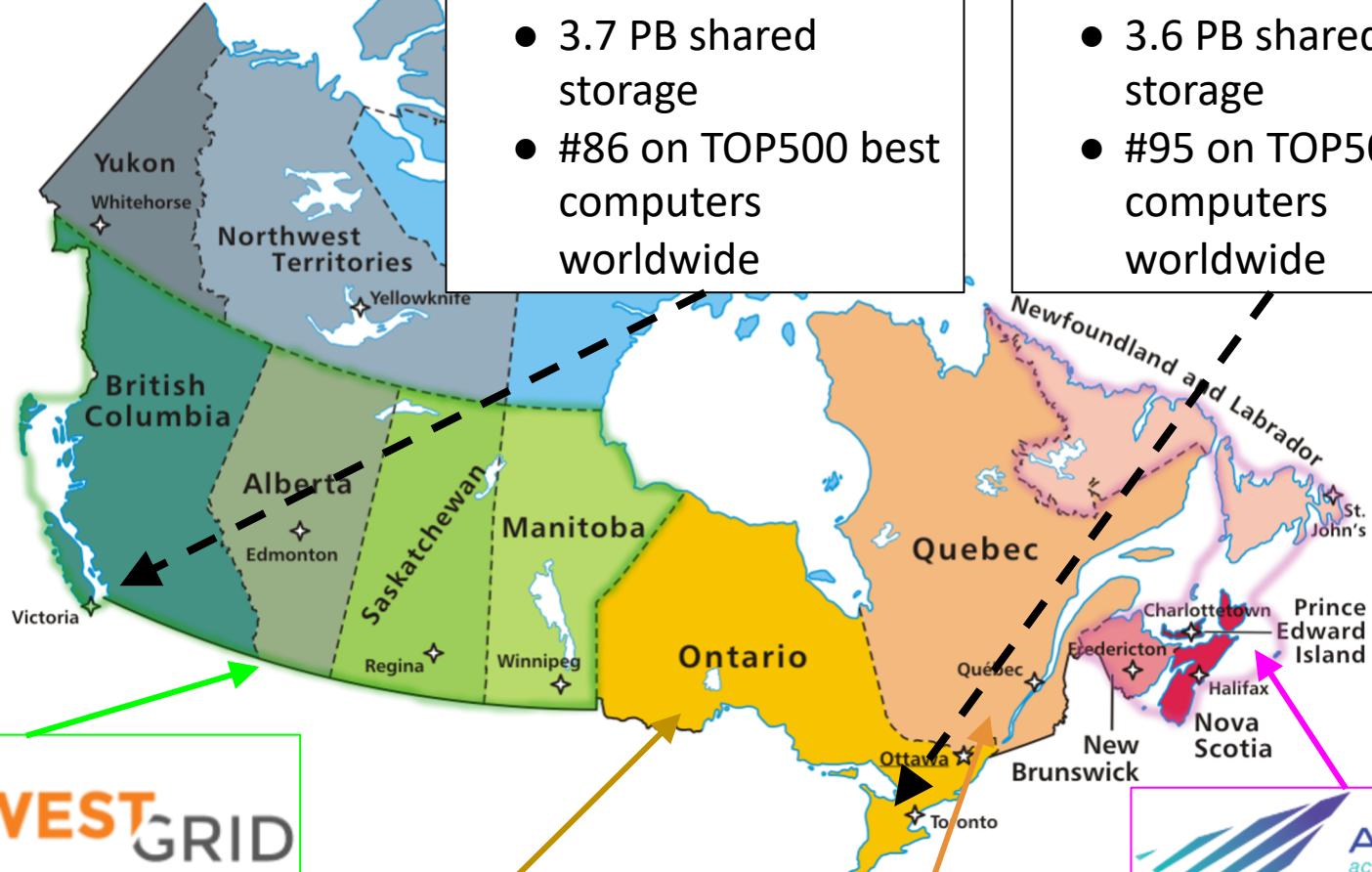
- **10** servers
- **350 TB** of storage
- **600** cpu cores
- **3 TB** of RAM

## Dedicated servers

- **12** Elasticsearch servers for big data management
- **4** virtualisation servers



compute | calcul  
canada | canada



## Cedar supercomputer

University Simon Fraser

- 59,776 shared CPU
- 3.7 PB shared storage
- #86 on TOP500 best computers worldwide

## Graham supercomputer

University of Waterloo

- 51,200 shared CPU
- 3.6 PB shared storage
- #95 on TOP500 best computers worldwide

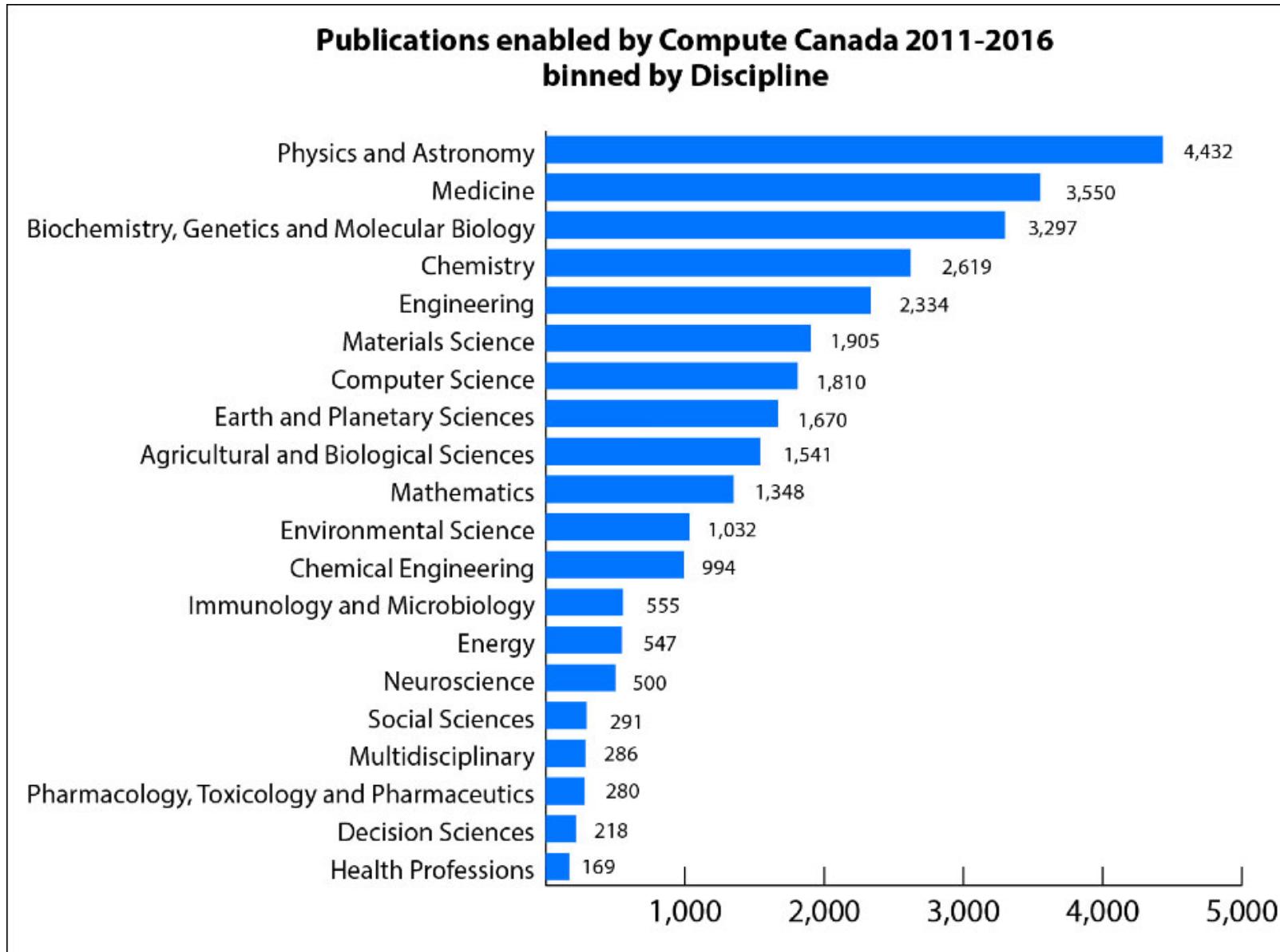


ACENET  
accelerate discovery

# Compute Canada allocation

- **78 cpu-year** in the beluga-compute cluster
- **300 TB** of PROJECT storage space in the ndc-calculquebec cluster
- **900 TB** of NEARLINE storage space in the ndc-calculquebec cluster
- **98 years-VCPU** in the arbutus-persistent-cloud cluster
- **10 TB** of INFONUAGIC storage space in the arbutus-persistent-cloud cluster
- **375 GB** of RAM in the arbutus-persistent-cloud cluster
- **11** volumes in the arbutus-persistent-cloud cluster
- **2** floating IP addresses in the arbutus-persistent-cloud cluster

# Publications enabled by Compute Canada

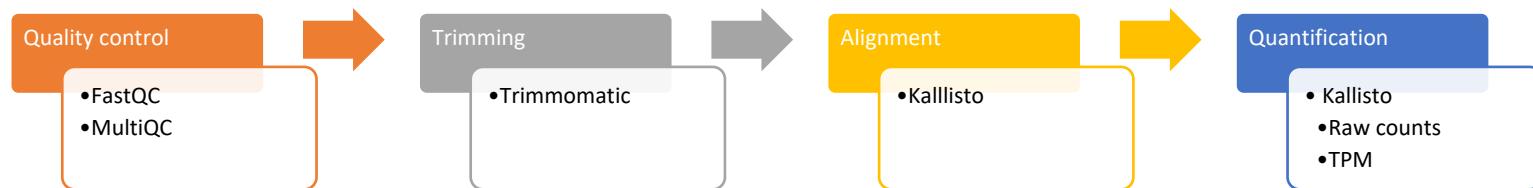


# Transcriptomics analyses from the Next Generation Sequencing Lab

# Transcriptomics/miRomics pipelines

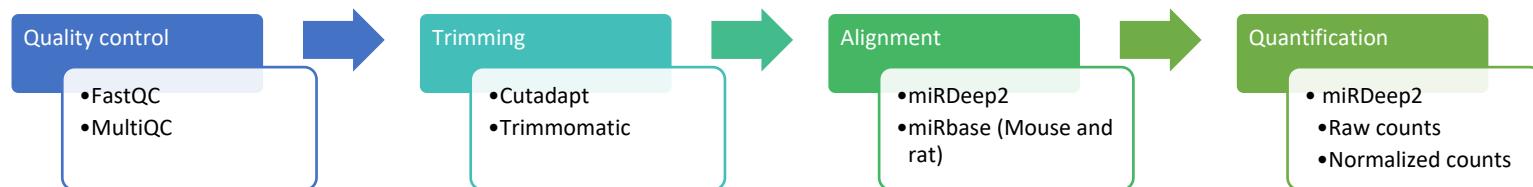
## 1. mRNA quantification

First, a quality control was performed using FastQC. Then, sequence reads were trimmed using Trimmomatic and aligned using Kallisto, which also produced gene quantification.



## 2. miRNA quantification

A quality control was performed using FastQC. Then, sequence reads were trimmed using cutadapt and trimmomatic, and were aligned with miRDeep2 against mouse non-coding RNAs. Then, miRNAs were quantified by miRDeep2 using rat and mouse mature miRNAs and pre-miRNAs. Rat and mouse miRNAs are from very close species, then it may be useful to map on these two species.



# R packages

# Bioconductor (R) packages development

## ■ Metagene

Compare the behavior of DNA-interacting proteins

## ■ Imetagene

A graphical interface for the metagene package

## ■ ENCODEXplorer

Compilation of ENCODE metadata

## ■ similaRpeak

calculates metrics which assign a level of similarity between ChIP-Seq profiles

## ■ Nucleosim

Generate synthetic nucleosome maps

## ■ consensusSeeker

Detection of consensus regions inside a group of experiences using genomic positions and genomic ranges

## ■ RJMCMC

Estimation of nucleosome positions for genome-wide profiling

## ■ TimeOmics

Integrate multi-Omics longitudinal data measured on the same biological samples and select key temporal features

## ■ Gwena

Gene co-expression network analysis and explore the results in a single pipeline

[Home](#) » [Bioconductor 3.2](#) » [Software Packages](#) » [metagene](#)

## metagene

platforms all | downloads top 50% | posts 0 | in Bioc 1 year  
build ok | commits 2.83 | test coverage 96%

A package to produce metagene plots

Bioconductor version: Release (3.2)

This package produces metagene plots to compare the behavior of DNA-interacting proteins at selected groups of genes/features. Bam files are used to increase the resolution. Multiple combination of group of bam files and/or group of genomic regions can be compared in a single analysis. Bootstrapping analysis is used to compare the groups and locate regions with statistically different enrichment profiles.

Author: Charles Joly Beauparlant <charles.joly-beauparlant at crchul.ulaval.ca>, Fabien Claude Lamaze <fabien.lamaze.1 at ulaval.ca>, Rawane Samb <rawane.samb.1 at ulaval.ca>, Astrid Louise Deschenes <Astrid-Louise.Deschenes at crchudequebec.ulaval.ca> and Arnaud Droit <arnaud.droit at crchul.ulaval.ca>.

Maintainer: Charles Joly Beauparlant <charles.joly-beauparlant at crchul.ulaval.ca>

Citation (from within R, enter `citation("metagene")`):

Beauparlant CJ, Lamaze FC, Samb R, Deschenes AL and Droit A (2014). *metagene: A package to produce metagene plots*. R package version 2.2.0.

## Installation

To install this package, start R and enter:

```
## try http:// if https:// URLs are not supported
source("https://bioconductor.org/biocLite.R")
biocLite("metagene")
```

## Documentation

To view documentation for the version of this package installed in your system, start R and enter:

```
browseVignettes("metagene")
```

[HTML](#)

metagene: a package to produce metagene plots

[PDF](#)

Reference Manual

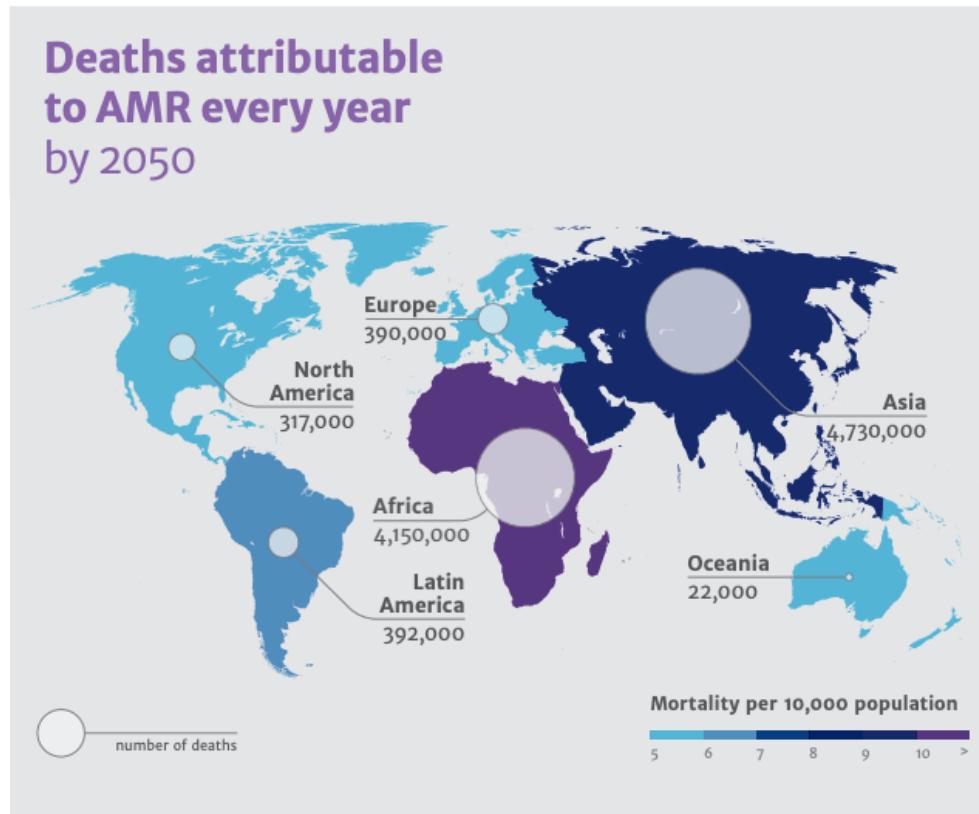
[Text](#)

NEWS

[Text](#)

LICENSE

# Antibiotics resistance



## GLOBAL

A failure to address the problem of antibiotic resistance could result in:



**10m deaths per year by 2050**

**Costing \$100 trillion in economic output**

## The next global challenge

# The causes of antibiotic resistance

**STOP OVERUSE AND  
MISUSE OF ANTIBIOTICS  
COMBAT RESISTANCE**



## CAUSES OF ANTIBIOTIC RESISTANCE



Antibiotic resistance happens when bacteria change and become resistant to the antibiotics used to treat the infections they cause.



Over-prescribing  
of antibiotics



Patients not finishing  
their treatment



Over-use of antibiotics in  
livestock and fish farming



Poor infection control  
in hospitals and clinics



Lack of hygiene and poor  
sanitation



Lack of new antibiotics  
being developed

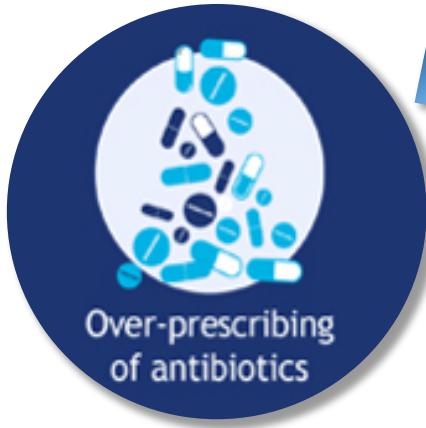
[www.who.int/drugresistance](http://www.who.int/drugresistance)

#AntibioticResistance

World Health Organization

# The causes of antibiotic resistance

**STOP OVERUSE AND  
MISUSE OF ANTIBIOTICS  
COMBAT RESISTANCE**



Over-prescribing  
of antibiotics

## CAUSES OF ANTIBIOTIC RESISTANCE

Antibiotic resistance happens when bacteria change and become resistant to the antibiotics used to treat the infections they cause.

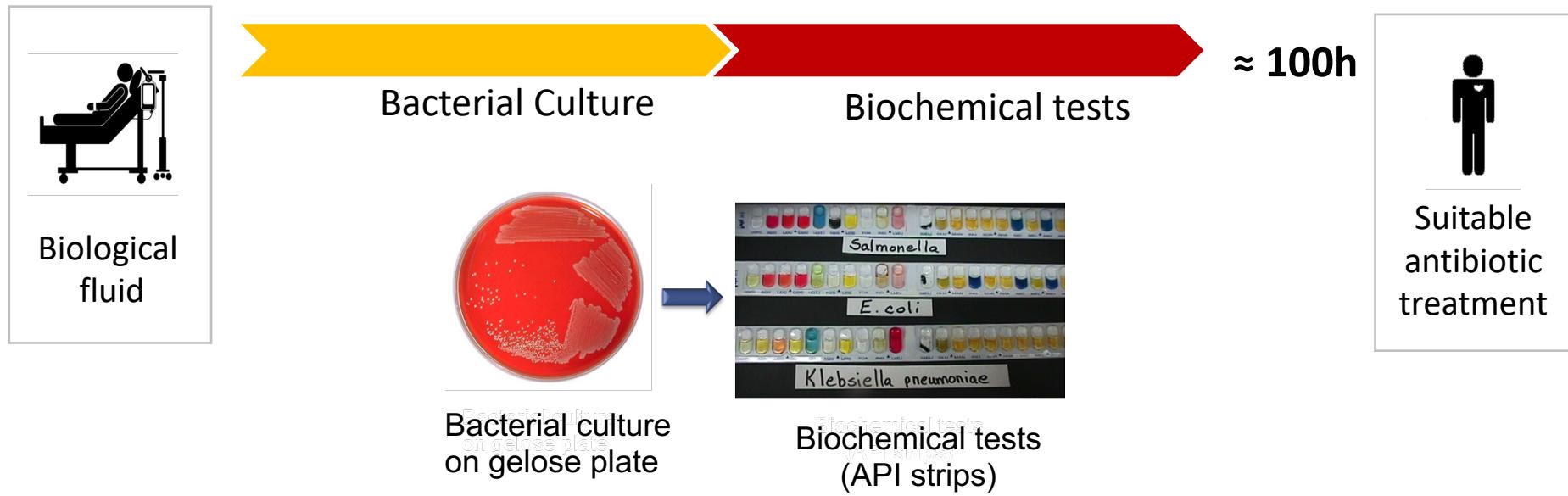


This diagram illustrates the causes of antibiotic resistance. It features a central title "CAUSES OF ANTIBIOTIC RESISTANCE" with a circular "HANDLE ANTIBIOTICS WITH CARE" logo. Below the title, a definition of antibiotic resistance is provided. Six circular icons, each with a corresponding caption, are arranged in a grid. A large blue arrow points from the "Over-prescribing of antibiotics" icon on the left towards the first icon in the grid. The icons and their captions are:

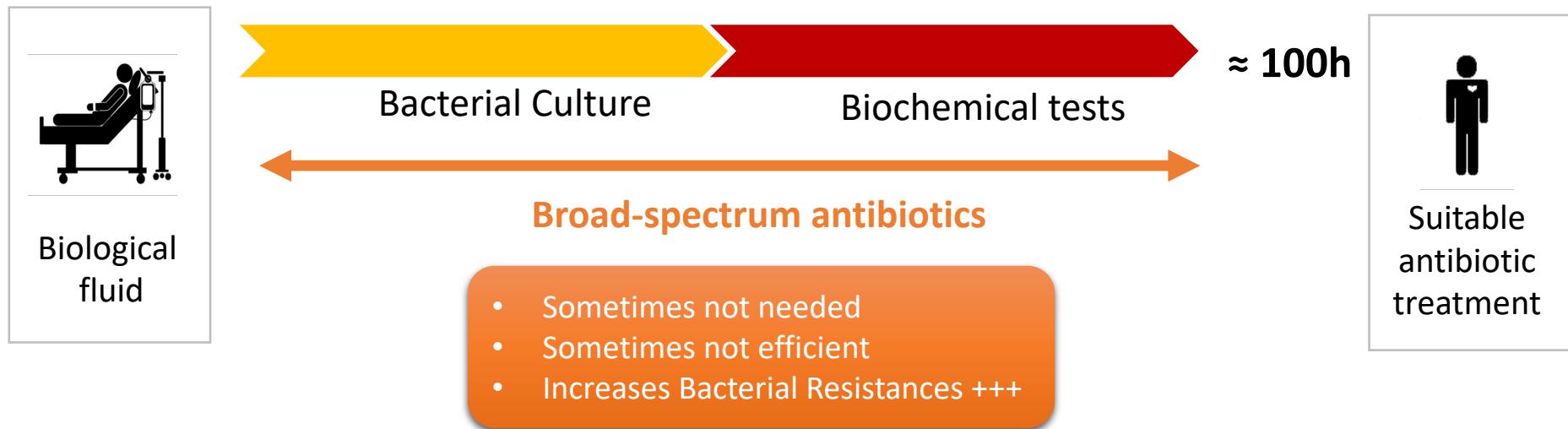
- Over-prescribing of antibiotics (blue pills icon)
- Poor infection control in hospitals and clinics (hospital room icon)
- Lack of hygiene and poor sanitation (handwashing icon)
- Lack of new antibiotics being developed (laboratory icon)
- Patients not finishing their treatment (pill bottle icon)
- Over-use of antibiotics in livestock and fish farming (cow and fish icon)

[www.who.int/drugresistance](http://www.who.int/drugresistance)  
**#AntibioticResistance**

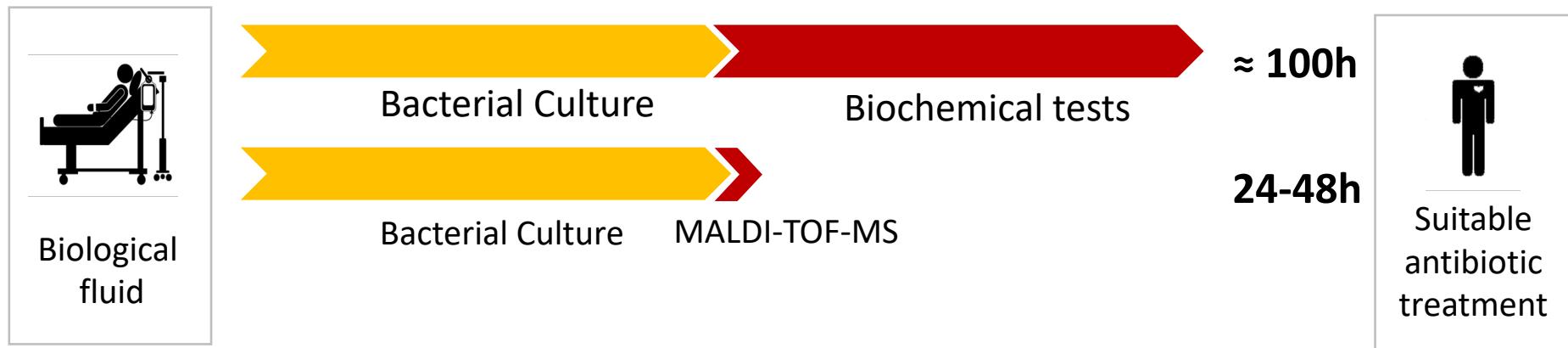
# Bacterial identification in clinics



# Bacterial identification in clinics



# Bacterial identification in clinics

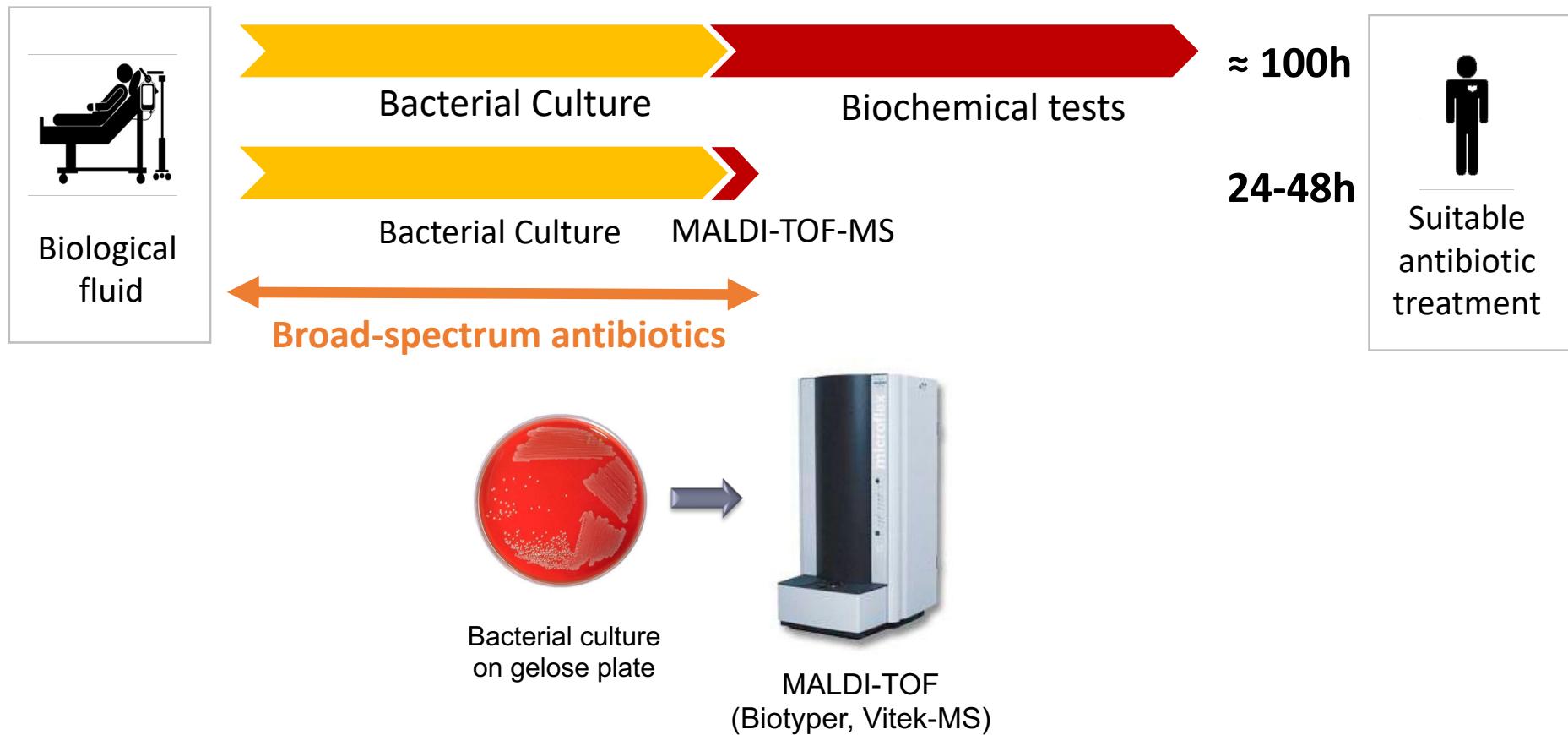


Bacterial culture  
on gelose plate

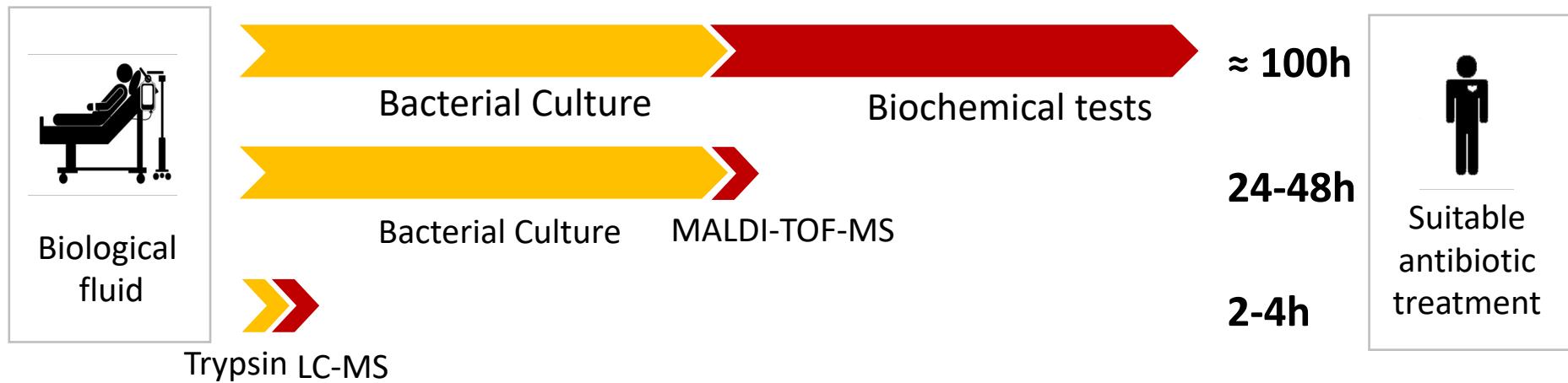


MALDI-TOF  
(Biotyper, Vitek-MS)

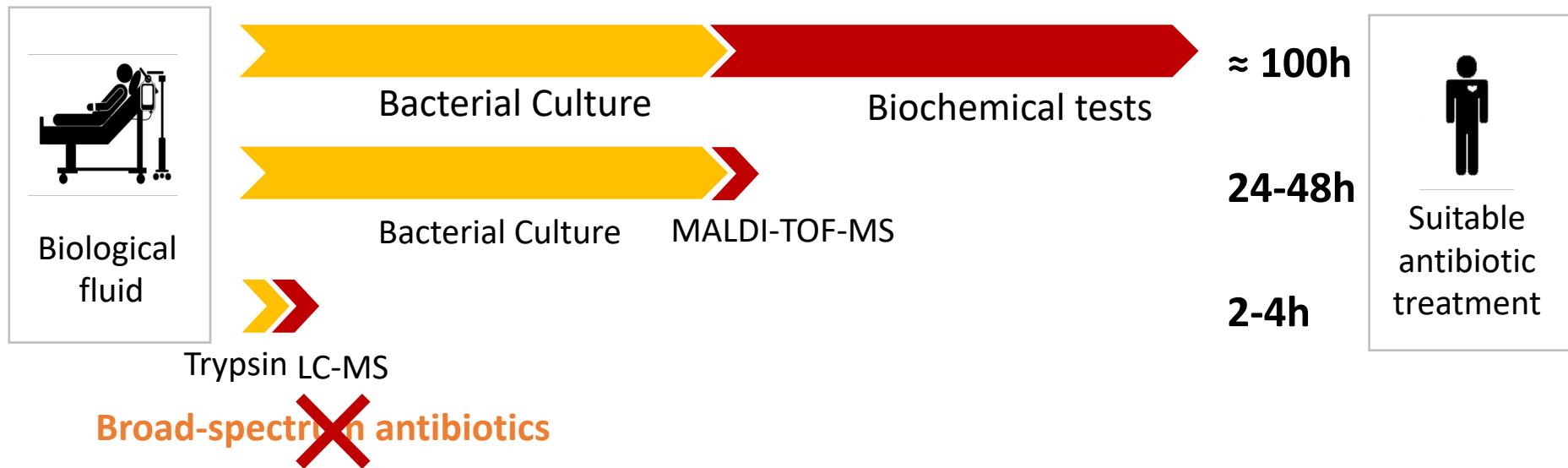
# Bacterial identification in clinics



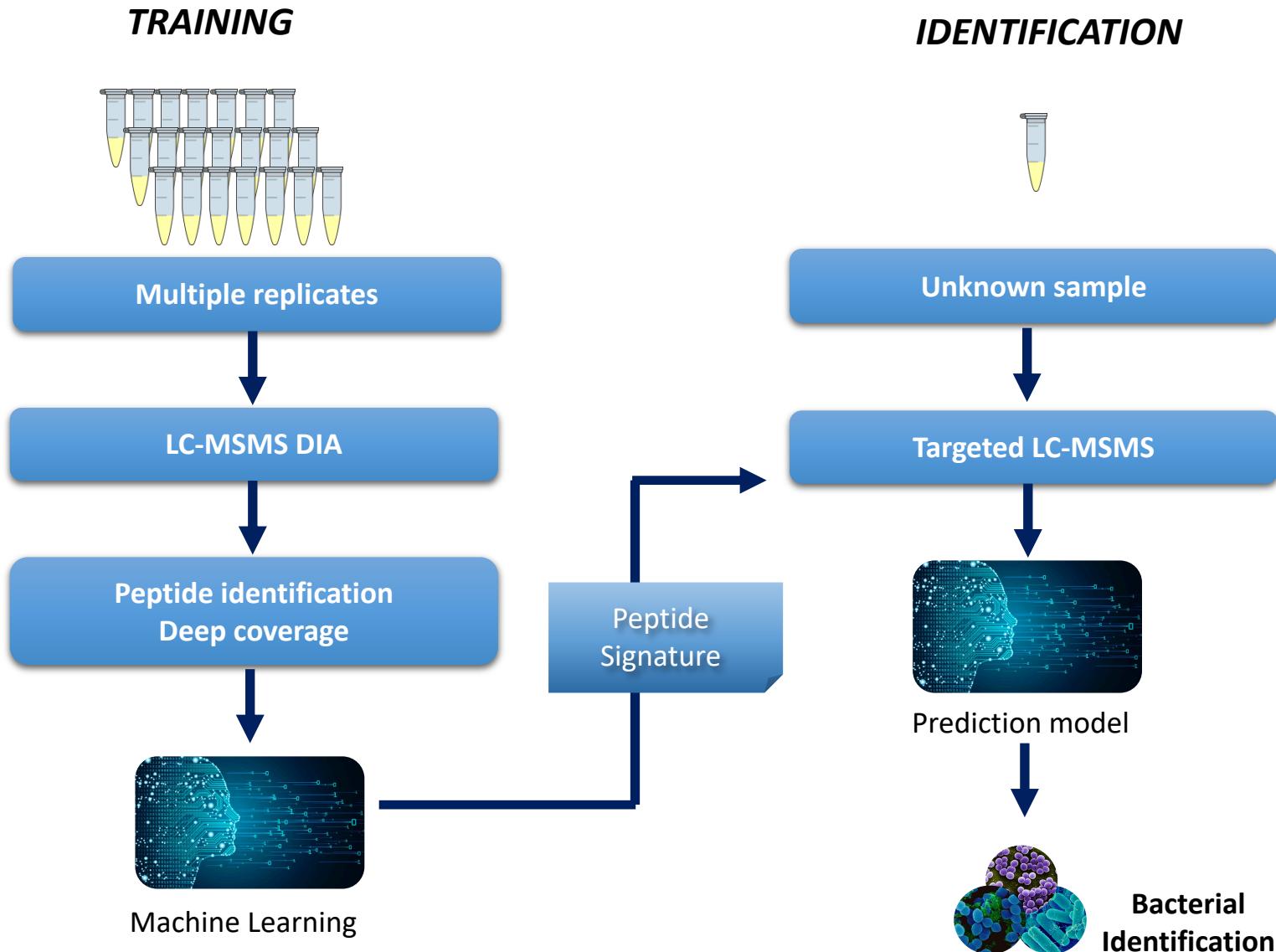
# Bacterial identification in clinics



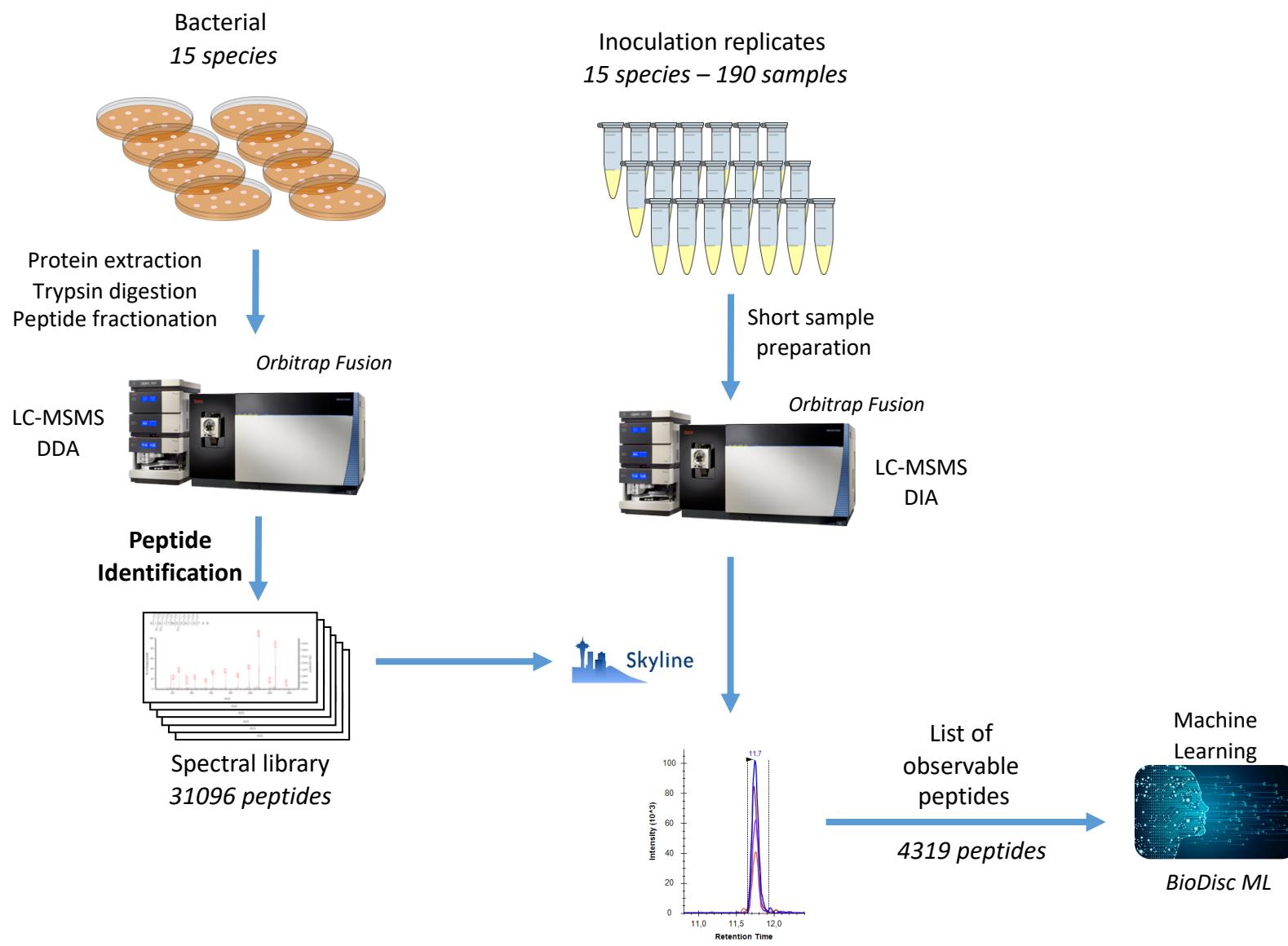
# Bacterial identification in clinics



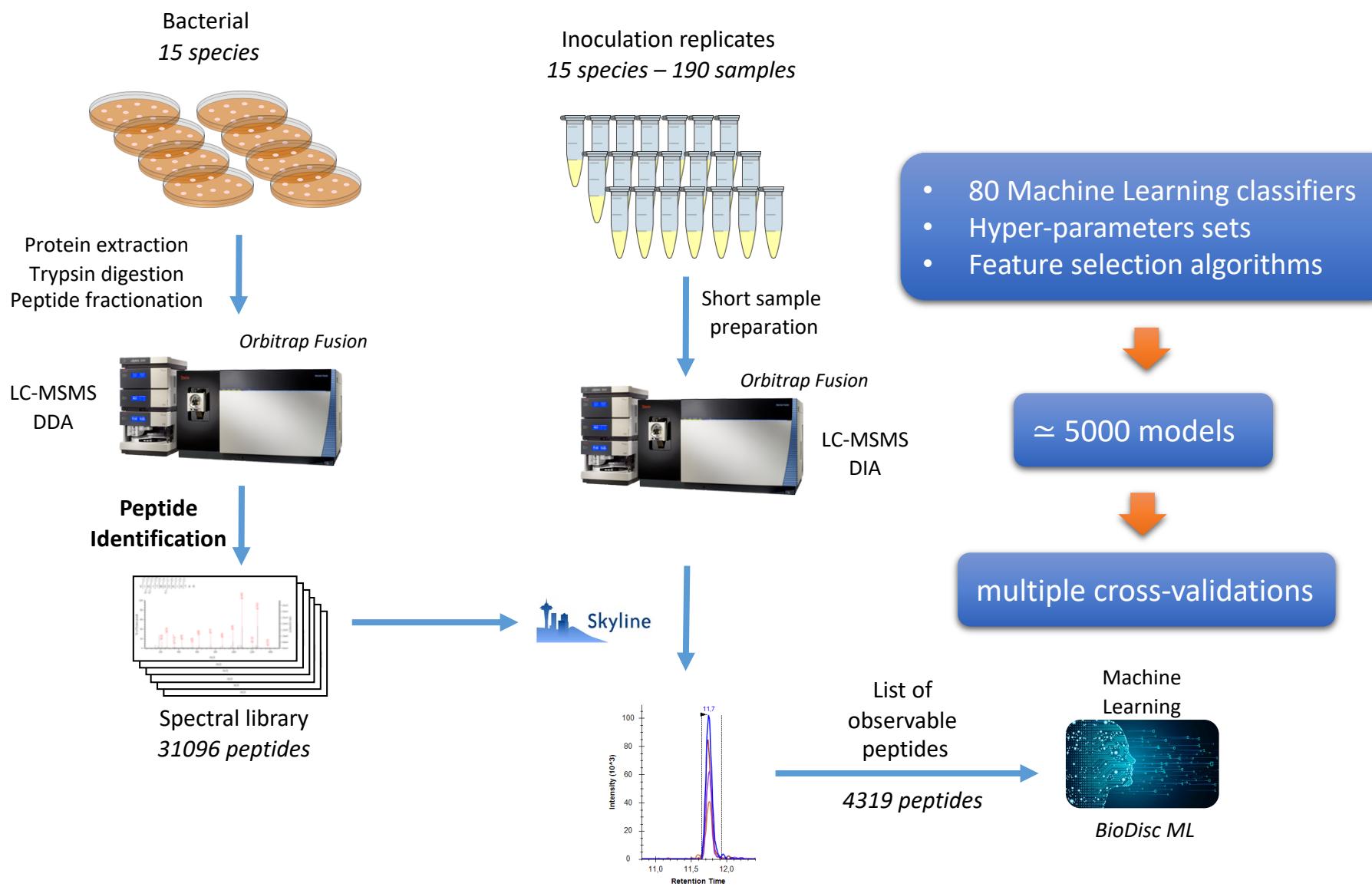
# Bacterial Identification by LC-MSMS + ML



# TRAINING STEP : Deep Proteome Coverage



# TRAINING STEP : Deep Proteome Coverage



# TRAINING STEP : Peptide signature

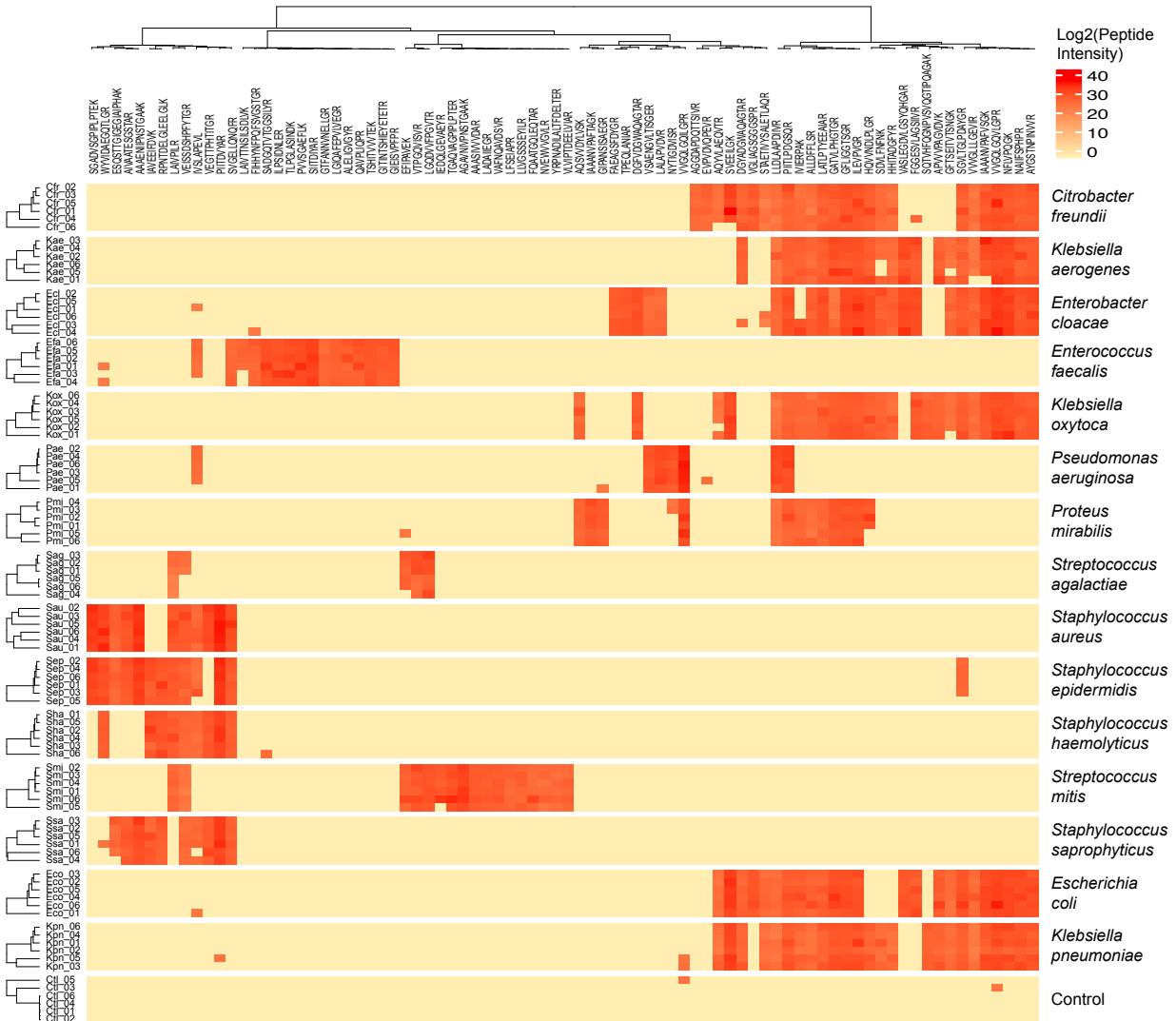


BioDiscML



82 Peptides  
Signature

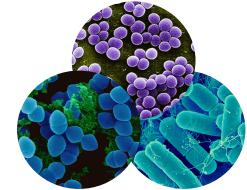
- Peptides shared by several species
- Each species has its own pattern



# Conclusion and perspectives



- LC-MS and Machine Learning allow the identification of bacteria in UTI in a short time without culture
- The method can work on triple quadrupole instruments
  - Development of quantification method for UTI on TSQ Altis
  - Development of new signatures (blood, milk, water....)
  - Detection of bacteria having specific resistances or virulences
- Improvement of the strategy using crude DIA signal and Machine Learning





# Assemblage de génome du caribou des bois (*Rangifer tarandus caribou*)



# Le Caribou

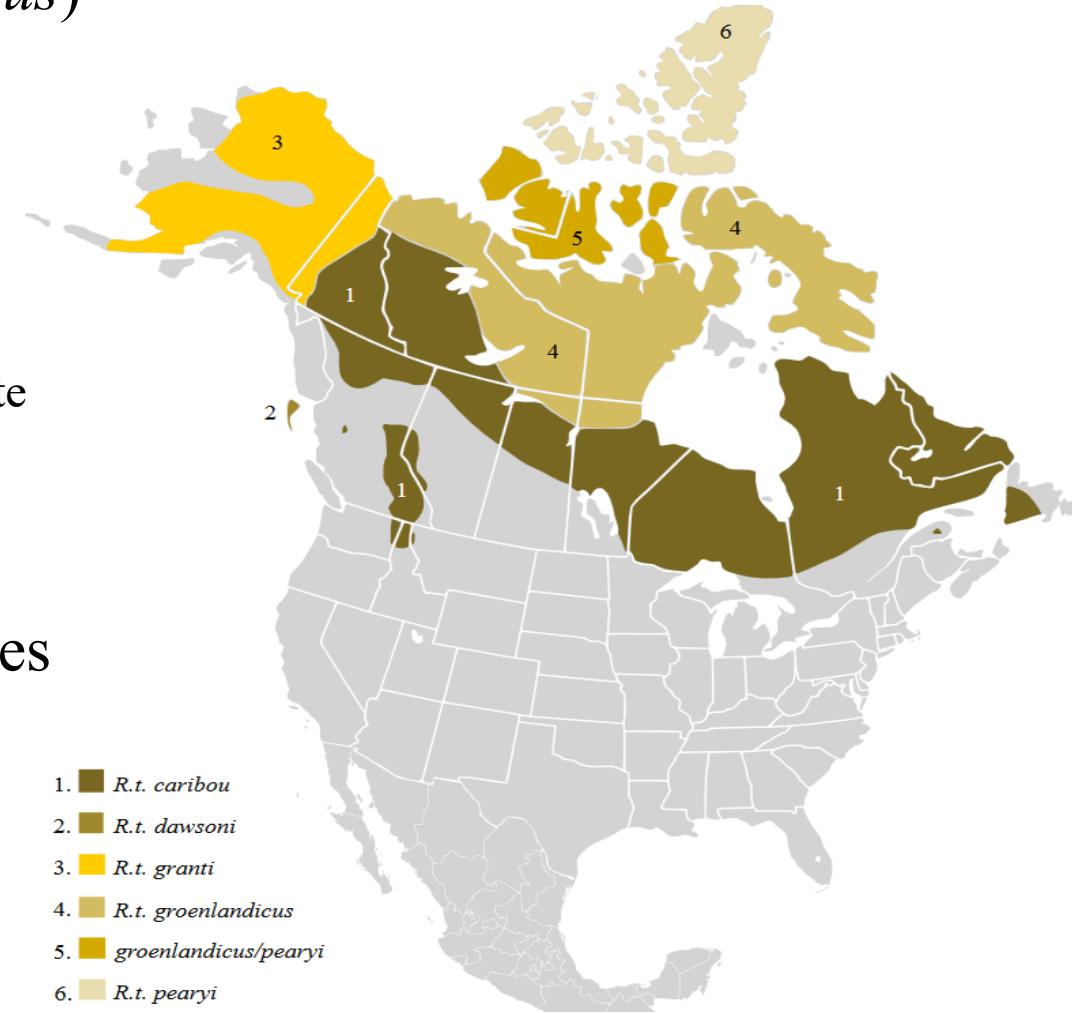
- Le caribou (*Rangifer tarandus*)

- Plusieurs ss-espèces:

- De Grant
- De la toundra
- Des îles du Nunavut et des TNO
- De l'archipel de la Reine-Charlotte (éteint)
- Des bois

- Plusieurs écotypes au sein des “bois”

- Migrateur (ou toundrique)
- Forestier
- Montagnard



# Alimentation et migration

- Alimentation: lichen dans des forêts de faible densité

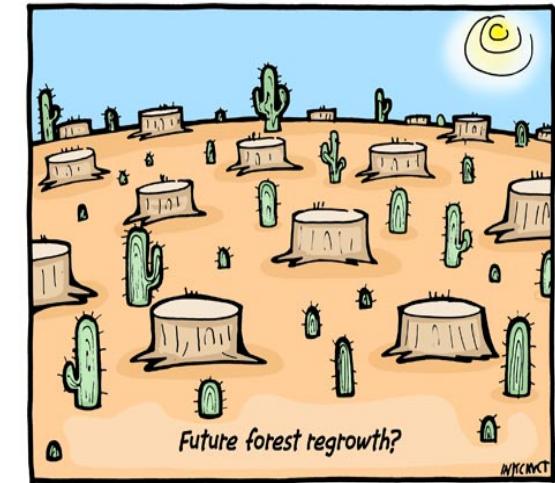
- Migration:

- Hiver en taiga (forêt)
- Été en toundra
- Deux grandes hardes au Québec

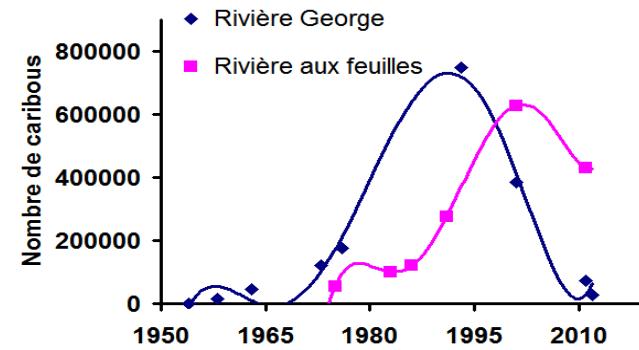


# Menaces

- Exploitation minière et forestière

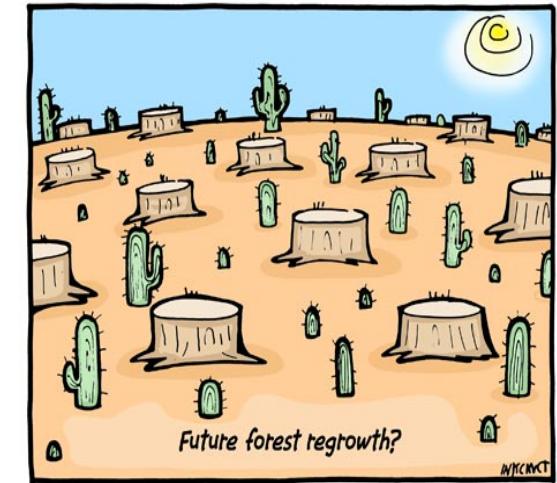


- Changements climatiques
- Baisse des tailles de populations



# Menaces

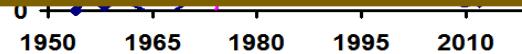
- Exploitation minière et forestière



- Changements climatiques

Nécessités de protection:

- > chasse interdite au Québec depuis 01/02/2018
- > exception: hardes migratrices par les 1<sup>ères</sup> nations



# Objectifs du projet

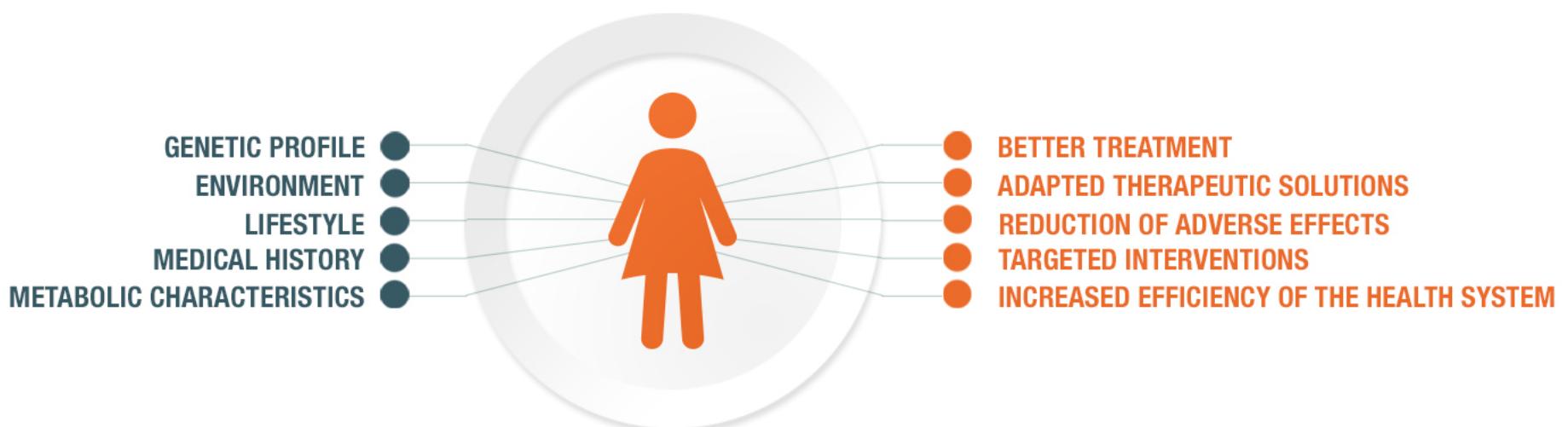
- Designer une biopuce à SNPs
  - Vision à long terme => polyvalence
  - Besoin d'un génome de haute qualité
- Définir les populations au niveau génétique (structure hiérarchique de populations)
- Mettre au point un outil de forensic pour déterminer la provenance d'un échantillon en cas de procès pour braconnage.



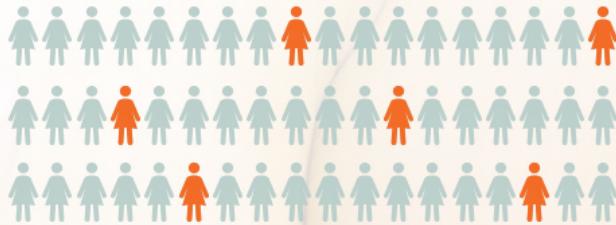
# Personalized medicine Breast Cancer

# PERSONALISED HEALTHCARE

- PERSONALISED
- PREDICTIVE
- PREVENTIVE
- PARTICIPATIVE



# BREAST CANCER BURDEN IN CANADA



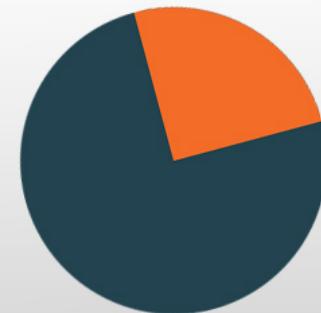
1/9 WOMEN WILL DEVELOP  
**BREAST CANCER**  
DURING HER LIFE

**5 000**

DEATHS ARE ATTRIBUTED TO  
THIS CANCER EACH YEAR

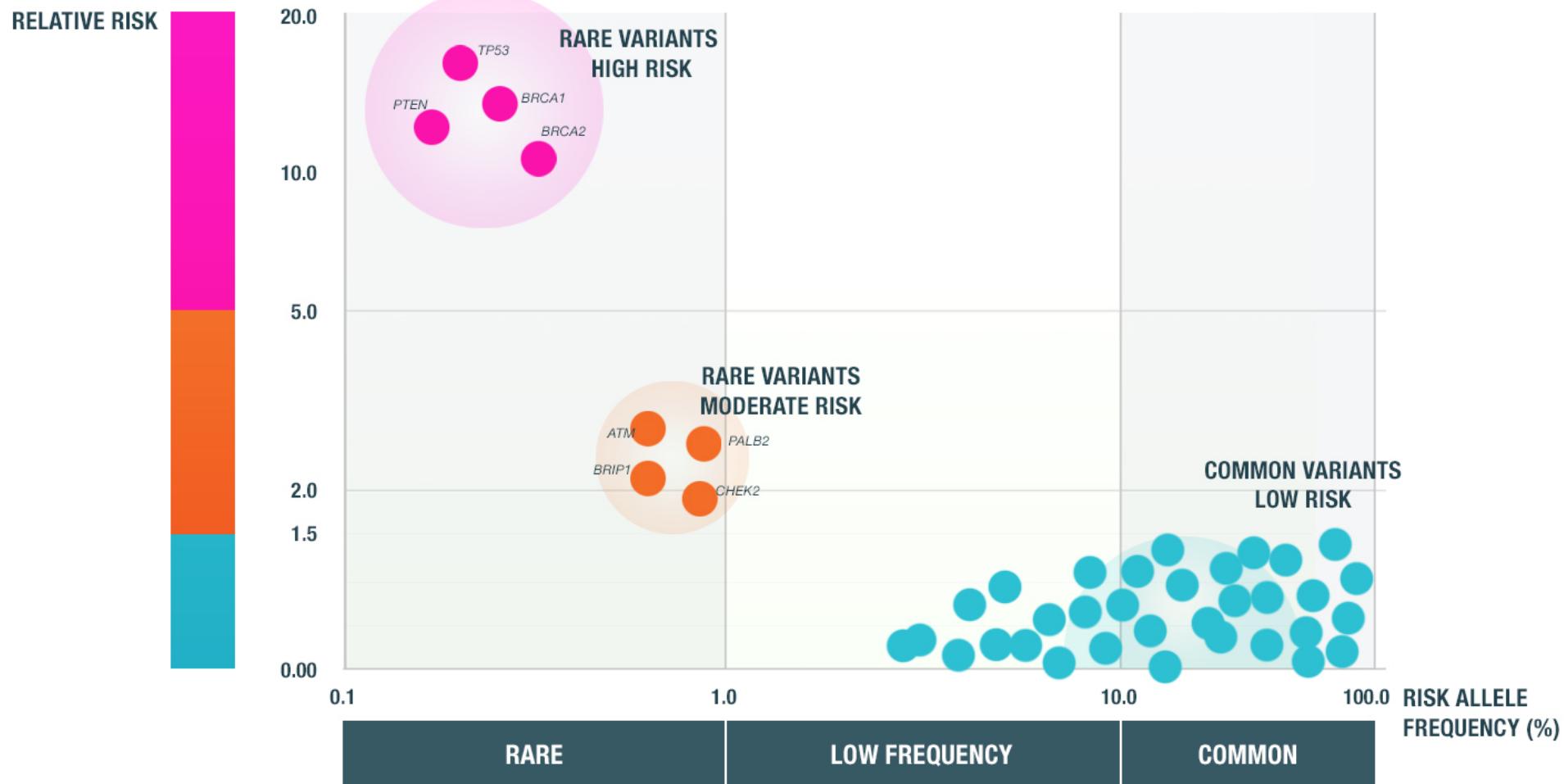
**1/5 CASE**

OCCUR IN WOMEN < 50 YEARS OLD



**36%** OF ALL CANCERS DIAGNOSED IN WOMEN  
BETWEEN THE AGES OF 30-49

# GENETIC VARIATIONS ASSOCIATED WITH BREAST CANCER RISK



# GENETIC RISK PROFILE

INTEGRATION OF  
DATA



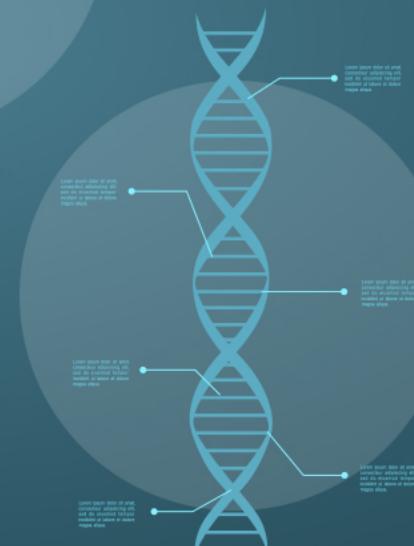
HUNDRED OF RESEARCH TEAMS  
FROM

**40**  
**COUNTRIES**



**150 000**  
PARTICIPANTS

**“ONCOARRAY”**



ANALYSIS OF  
GENOMIC  
DATA



MEDICAL HISTORY  
TUMOR PATHOLOGY  
TREATMENT RESPONSE  
ENVIRONMENT  
FAMILY HISTORY

# Computational genomics challenges

Royaume-Uni

Pays-Bas

Allemagne

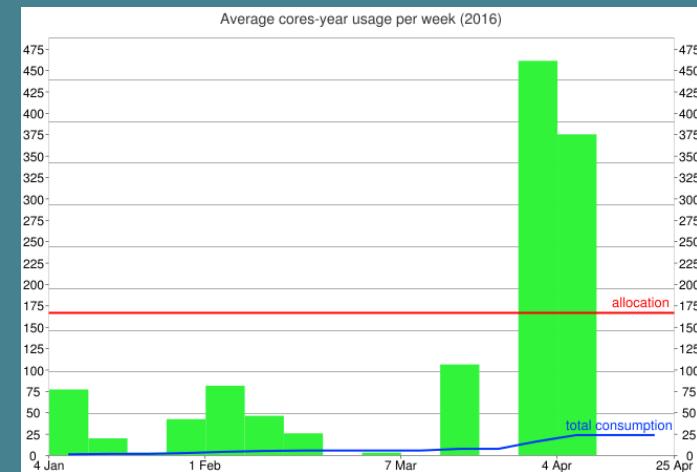
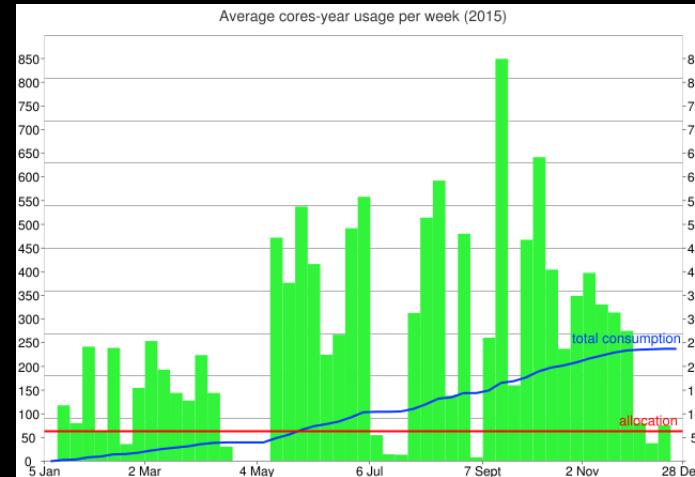
Canada (ON, QC)

États-Unis (UT)

Projet international

5 050 échantillons

Séquençage Massivement  
Parallèle



Consommation 2015 - mi 2016 : ~300 coeurs-année  
= 150 ans de calcul  
sur un ordinateur personnel moyen

# Timeline Perspective validation

Sample collection

Selection 251 genes

Targeted re-sequencing & Statistical analysis

May-July  
2016

September  
2016

September  
2017

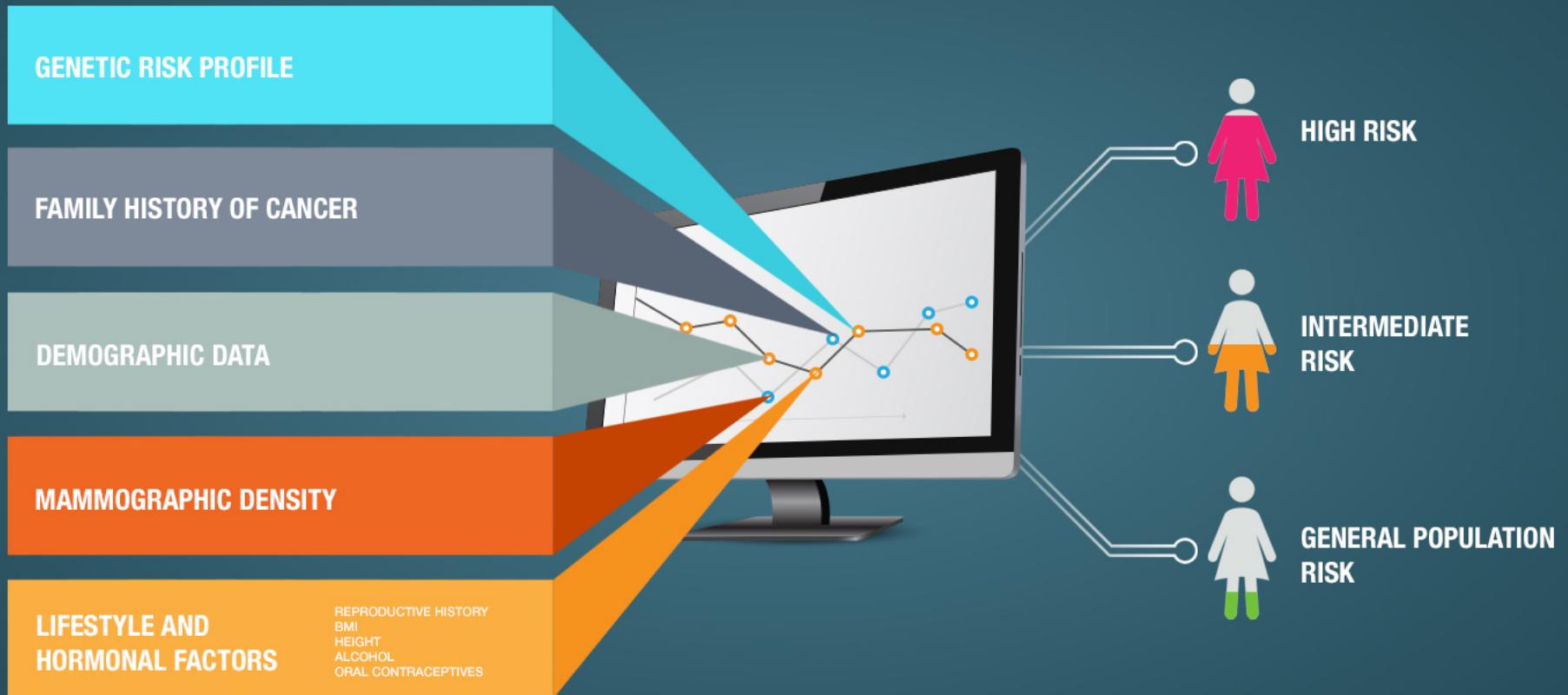


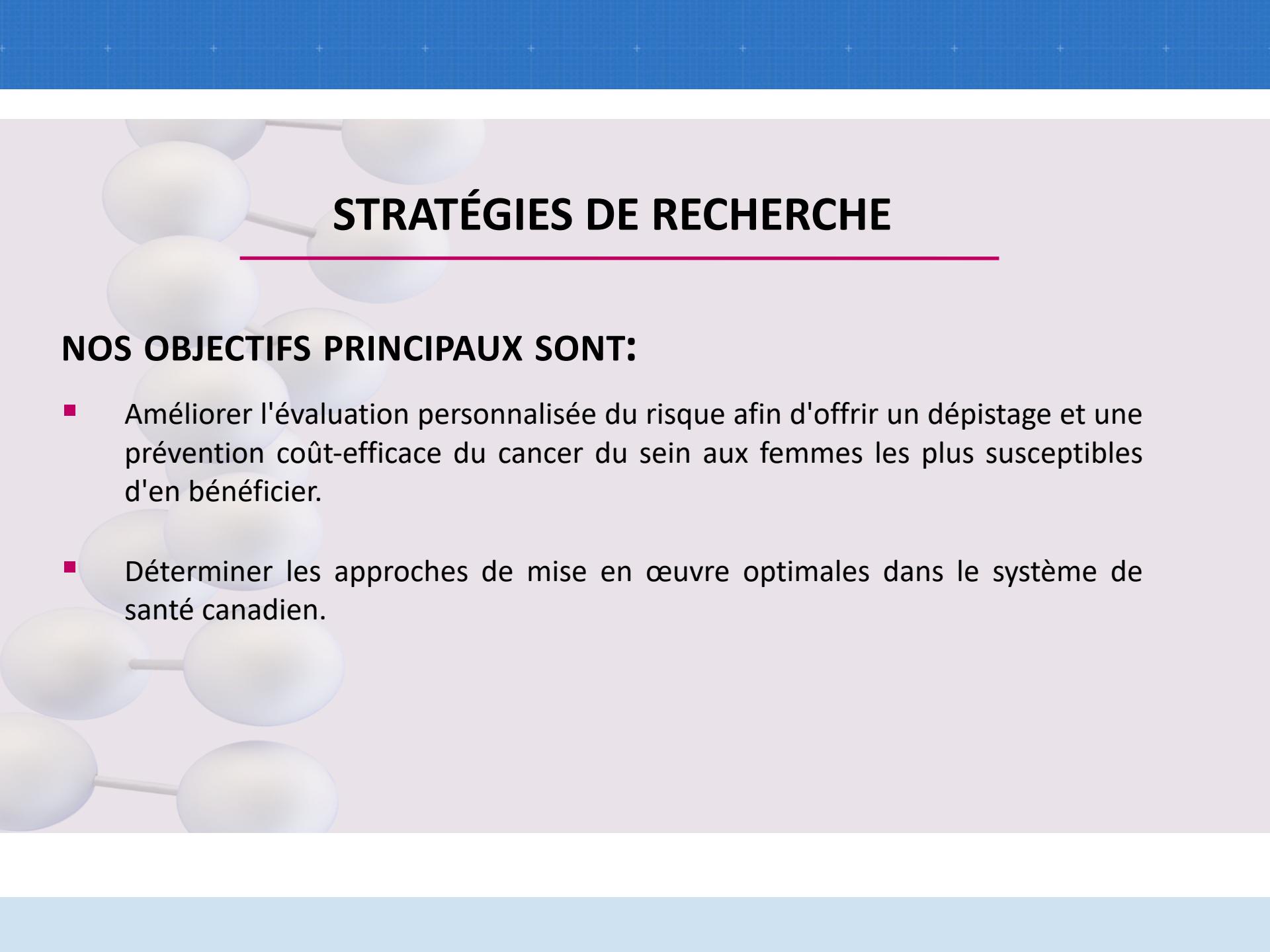
# GO Results - 251 genes

Y a-t-il des fonctions, cible d'intérêt dans le cancer ?

GO biological process complete	#	upload 1 (▼ Hierarchy NEW! ?)				
		#	expected	Fold Enrichment	+/-	P value
<a href="#">strand displacement</a>	<a href="#">26</a>	<a href="#">13</a>	.33	39.13	+	4.45E-13
↳ <a href="#">DNA recombination</a>	<a href="#">226</a>	<a href="#">24</a>	2.89	8.31	+	4.17E-11
↳ <a href="#">DNA metabolic process</a>	<a href="#">787</a>	<a href="#">55</a>	10.06	5.47	+	7.83E-21
↳ <a href="#">cellular macromolecule metabolic process</a>	<a href="#">6789</a>	<a href="#">124</a>	86.76	1.43	+	1.25E-02
↳ <a href="#">macromolecule metabolic process</a>	<a href="#">7494</a>	<a href="#">133</a>	95.77	1.39	+	1.83E-02
↳ <a href="#">nucleic acid metabolic process</a>	<a href="#">3987</a>	<a href="#">82</a>	50.95	1.61	+	3.03E-02
↳ <a href="#">nucleobase-containing compound metabolic process</a>	<a href="#">4541</a>	<a href="#">93</a>	58.03	1.60	+	5.38E-03
↳ <a href="#">heterocycle metabolic process</a>	<a href="#">4678</a>	<a href="#">93</a>	59.78	1.56	+	2.05E-02
↳ <a href="#">cellular aromatic compound metabolic process</a>	<a href="#">4735</a>	<a href="#">94</a>	60.51	1.55	+	1.86E-02
<a href="#">DNA double-strand break processing</a>	<a href="#">19</a>	<a href="#">8</a>	.24	32.95	+	1.82E-06
↳ <a href="#">double-strand break repair</a>	<a href="#">161</a>	<a href="#">25</a>	2.06	12.15	+	1.90E-15
↳ <a href="#">DNA repair</a>	<a href="#">487</a>	<a href="#">48</a>	6.22	7.71	+	6.25E-24
↳ <a href="#">cellular response to DNA damage stimulus</a>	<a href="#">735</a>	<a href="#">54</a>	9.39	5.75	+	2.22E-21
↳ <a href="#">cellular response to stress</a>	<a href="#">1618</a>	<a href="#">76</a>	20.68	3.68	+	7.94E-20
↳ <a href="#">cellular response to stimulus</a>	<a href="#">6211</a>	<a href="#">144</a>	79.37	1.81	+	1.37E-12
↳ <a href="#">response to stimulus</a>	<a href="#">7683</a>	<a href="#">153</a>	98.18	1.56	+	6.97E-08
↳ <a href="#">response to stress</a>	<a href="#">3365</a>	<a href="#">98</a>	43.00	2.28	+	2.69E-12

# RISK PROFILING





## **STRATÉGIES DE RECHERCHE**

---

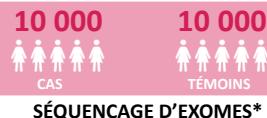
### **NOS OBJECTIFS PRINCIPAUX SONT:**

- Améliorer l'évaluation personnalisée du risque afin d'offrir un dépistage et une prévention coût-efficace du cancer du sein aux femmes les plus susceptibles d'en bénéficier.
- Déterminer les approches de mise en œuvre optimales dans le système de santé canadien.

# CLINIQUES D'ONCOGÉNÉTIQUE

## ACTIVITÉ 1

IDENTIFICATION DE NOUVEAUX GÈNES DE SUSCEPTIBILITÉ AU CANCER DU SEIN



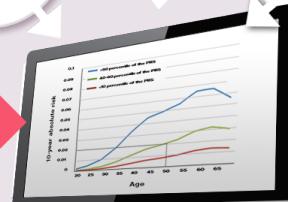
CRIBLAGE DE GÈNES CANDIDATS

ESTIMATIONS PRÉCISES DU RISQUE ASSOCIÉ AUX GÈNES DE SUSCEPTIBILITÉ AU CANCER DU SEIN

\* ÉVALUATION FONCTIONNELLE DE VARIANTES RARES

## ACTIVITÉ 2

SCORE DE RISQUE POLYGÉNIQUE AMÉLIORÉ  
ESTIMATIONS DU RISQUE CHEZ LES NON-EUROPEENS



TEST DE PANEL DE GÈNES

PREDICTION DU RISQUE SELON LE STATUT RH

PATHOGÉNICITÉ DES VARIANTES

NOUVEAUX GÈNES DE SUSCEPTIBILITÉ

ÉVALUATION DU RISQUE PERSONNALISÉ

EFFETS COMBINÉS DES VARIANTES COMMUNES ET RARES

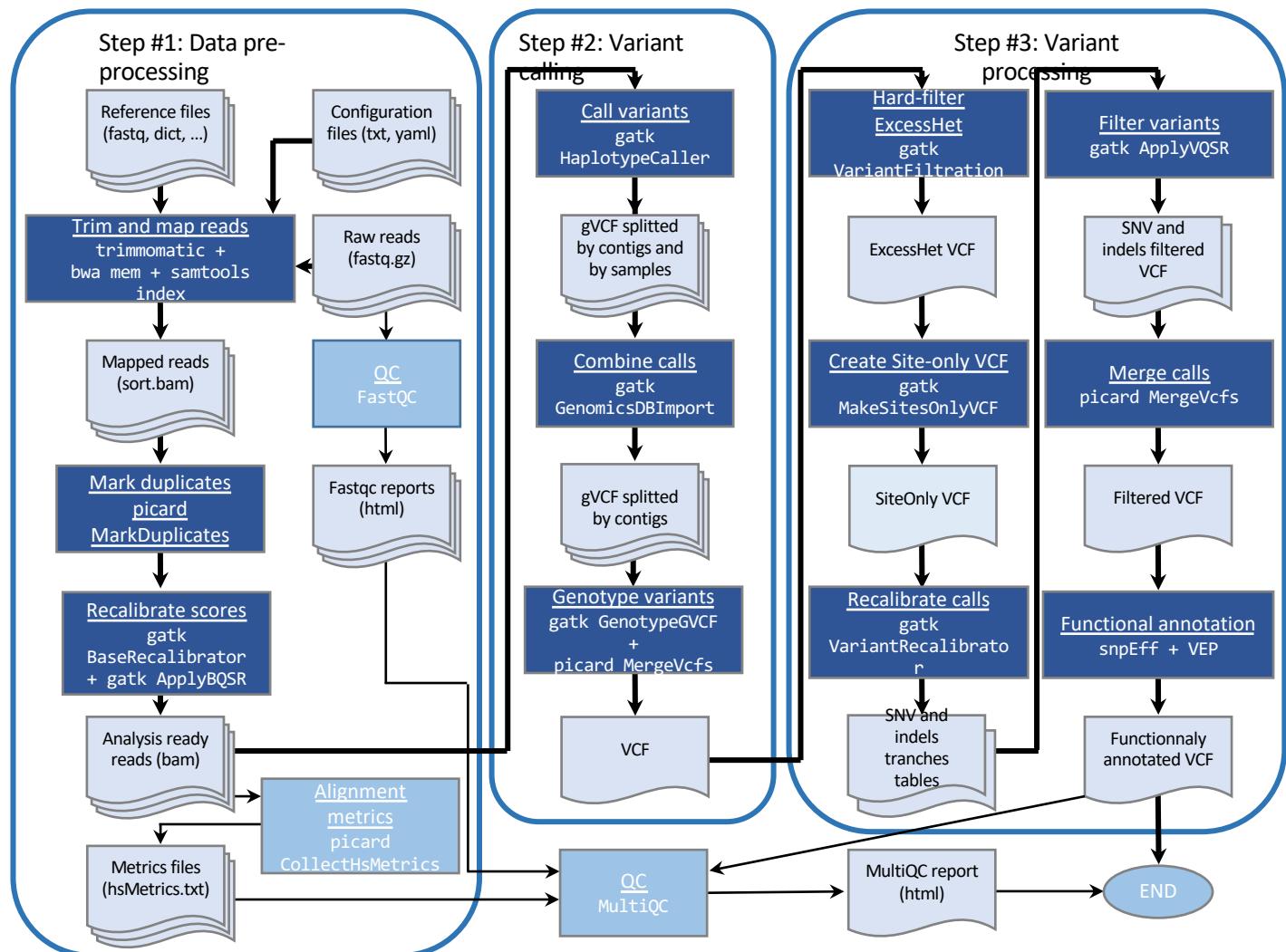
## PRÉDICTION DU RISQUE PLUS PRÉCISE



## AMÉLIORER:

- CONSEIL GÉNÉTIQUE
- PROCESSUS DE DÉCISION PARTAGÉE À PROPOS DU DÉPISTAGE ET DES STRATÉGIES DE RÉDUCTION DU RISQUE
- EFFICACITÉ DE LA PRÉVENTION PRIMAIRE

\* en collaboration avec le projet international BRIDGES



# Multi-omics analyses in Prostate Cancer

# Prostate cancer

- Prostate cancer affects about 1 in 7 men
- Cancer most common specific sex in men
- Patients are almost always brought to prostatectomy
- Biochemical recurrence may occur

*What therapeutics strategies ?*

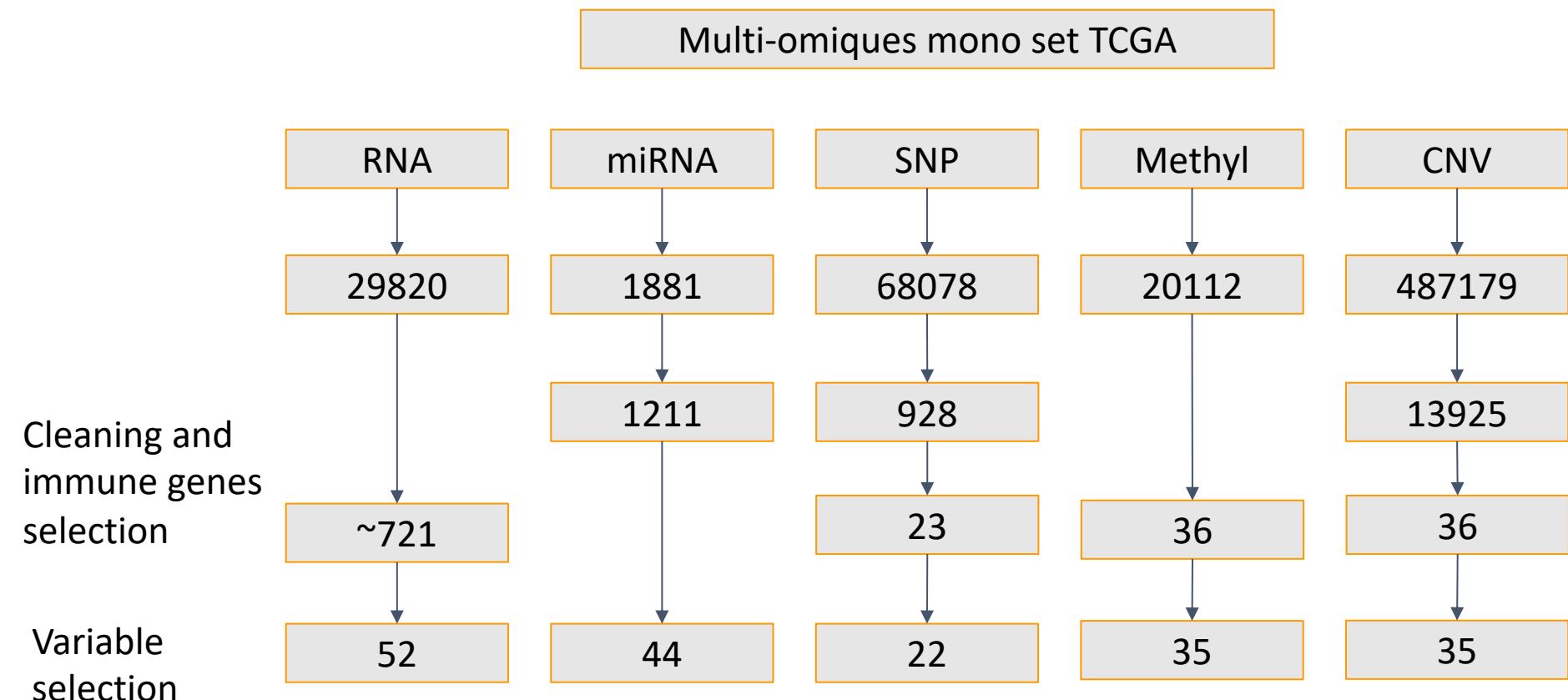
- prostatectomy
- Hormonotherapy
- immunotherapy

# Study objectives

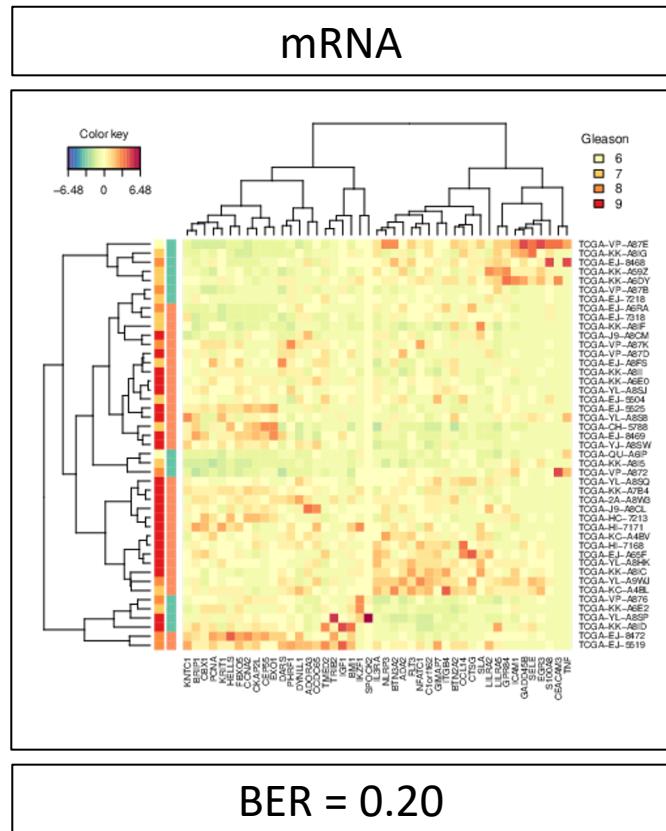
**Hypothesis:** The use of different omics datasets will allow to find recurring and significant biological markers related to the immunity to propose new immunological targets.

**Objectives:** 1) Retrieve datasets 2) Treat them uniformly 3) Target biomarkers of immunity 4) Choose appropriate statistical techniques

# Data curation and variable selection with MixOmics

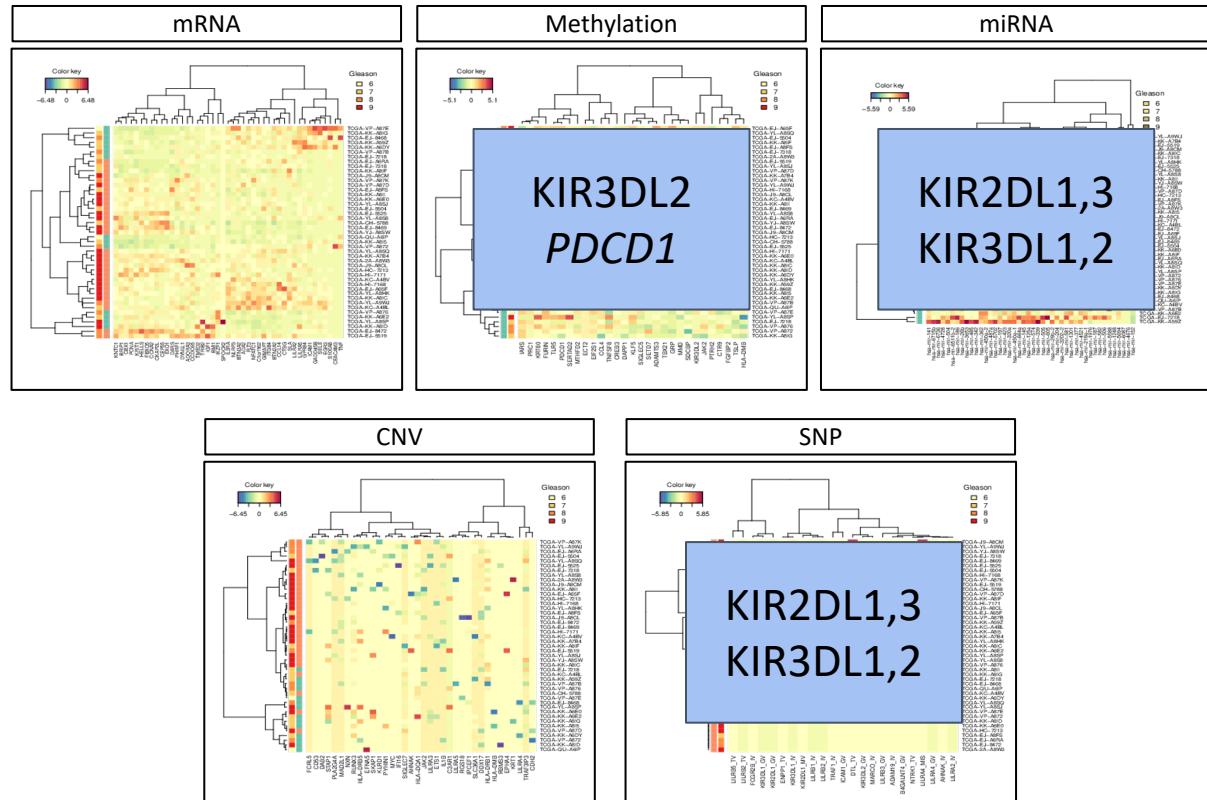


# mRNA signature



# Most interesting genes families

Gene	Type	Count
ADAM19	intron variant	2
AHNAK	intron variant	2
B4GALNT4	gene variant	3
DTL	transcript variant	2
ENPP1	transcript variant	2
FCGR2B	intron variant	2
ICAM1	gene variant	2
KIR2DL1	gene variant	4
KIR2DL1	missense variant	2
KIR2DL3	intron variant	10
KIR3DL1	gene variant	3
KIR3DL1	intron variant	19
KIR3DL2	gene variant	2
LILRA2	intron variant	3
LILRA4	gene variant	4
LILRA4	missense variant	2
LILRB1	intron variant	6
LILRB2	intron variant	7
LILRB2	transcript variant	3
LILRB3	gene variant	2
LILRB5	transcript variant	3
MARCO	intron variant	2
NTRK1	transcript variant	2
TRAF1	intron variant	2



# Development of Multi-OMICS strategies for the study of Human Gut Microbiota

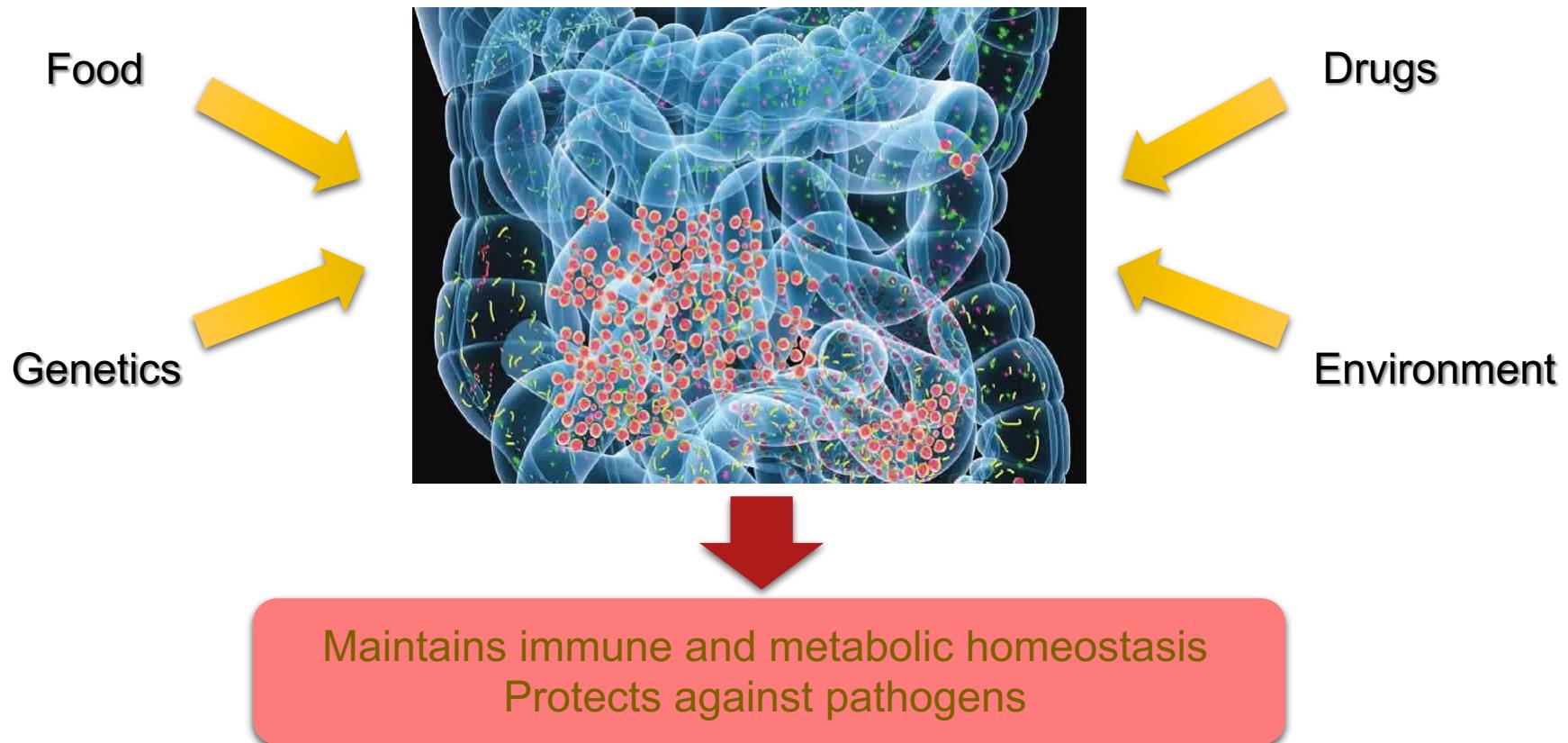


CENTRE DE RECHERCHE  
EN INFECTIOLOGIE

# Gut Microbiota

$10^{14}$  microorganisms (ratio 1:1 with human cells)

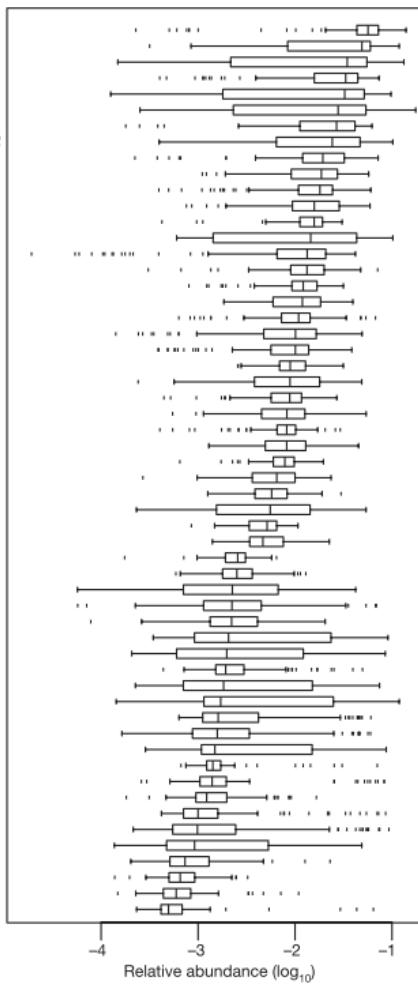
Hundreds/Thousands of species but 4 major phyla :  
93.5% belonged to Bacteroidetes, Firmicutes, Actinobacteria and Proteobacteria



Dysbiosis : Inflammatory diseases , diabete, Infection

# Model for studying the effect of antibiotics on gut microbiota

Bacteroides uniformis  
Alistipes putredinis  
Parabacteroides merdae  
Dorea longicatena  
Ruminococcus bromii L2-63  
Bacteroides caccae  
Clostridium sp. SS2-1  
Bacteroides thetaiotaomicron VPI-5482  
Eubacterium hallii  
Ruminococcus torques L2-14  
Unknown sp. SS3 4  
Ruminococcus sp. SR1 5  
Faecalibacterium prausnitzii SL3 3  
Ruminococcus lactaris  
Collinsella aerofaciens  
Dorea formicigenerans  
Bacteroides vulgaris ATCC 8482  
Roseburia intestinalis M50 1  
Bacteroides sp. 2\_1\_7  
Eubacterium siraeum 70 3  
Parabacteroides distasonis ATCC 8503  
Bacteroides sp. 9\_1\_42FAA  
Bacteroides ovatus  
Bacteroides sp. 4\_3\_47FAA  
Bacteroides sp. 2\_2\_4  
Eubacterium rectale M104 1  
Bacteroides xylinosolvens XB1A  
Coprococcus comes SL7 1  
Bacteroides sp. D1  
Bacteroides sp. D4  
Eubacterium ventriosum  
Bacteroides dorei  
Ruminococcus obaeum A2-162  
Subdoligranulum variabile  
Bacteroides capillosus  
Streptococcus thermophilus LMD-9  
Clostridium leptum  
Holdemania filiformis  
Bacteroides stercoris  
Coprococcus eutactus  
Clostridium sp. M62 1  
Bacteroides eggertii  
Butyrivibrio crossotus  
Bacteroides finegoldii  
Parabacteroides johnsonii  
Clostridium sp. L2-50  
Clostridium nexile  
Bacteroides pectiniphilus  
Anaerotruncus colihominis  
Ruminococcus gravus  
Bacteroides intestinalis  
Bacteroides fragilis 3\_1\_12  
Clostridium asparagiforme  
Enterococcus faecalis TX0104  
Clostridium scindens  
Blautia hansenii

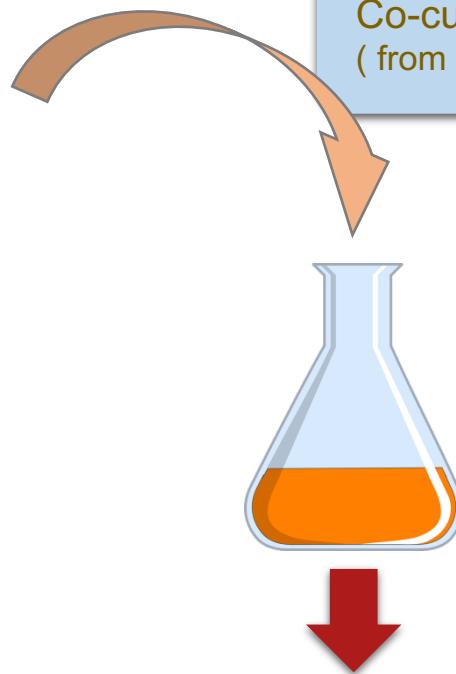


from Qin J. et al, 2010

Selection of species most frequently found and most representative of gut microbiota

## Artificial gut microbiota

Co-culture of selected species  
(from Institut Pasteur and CRI collections)

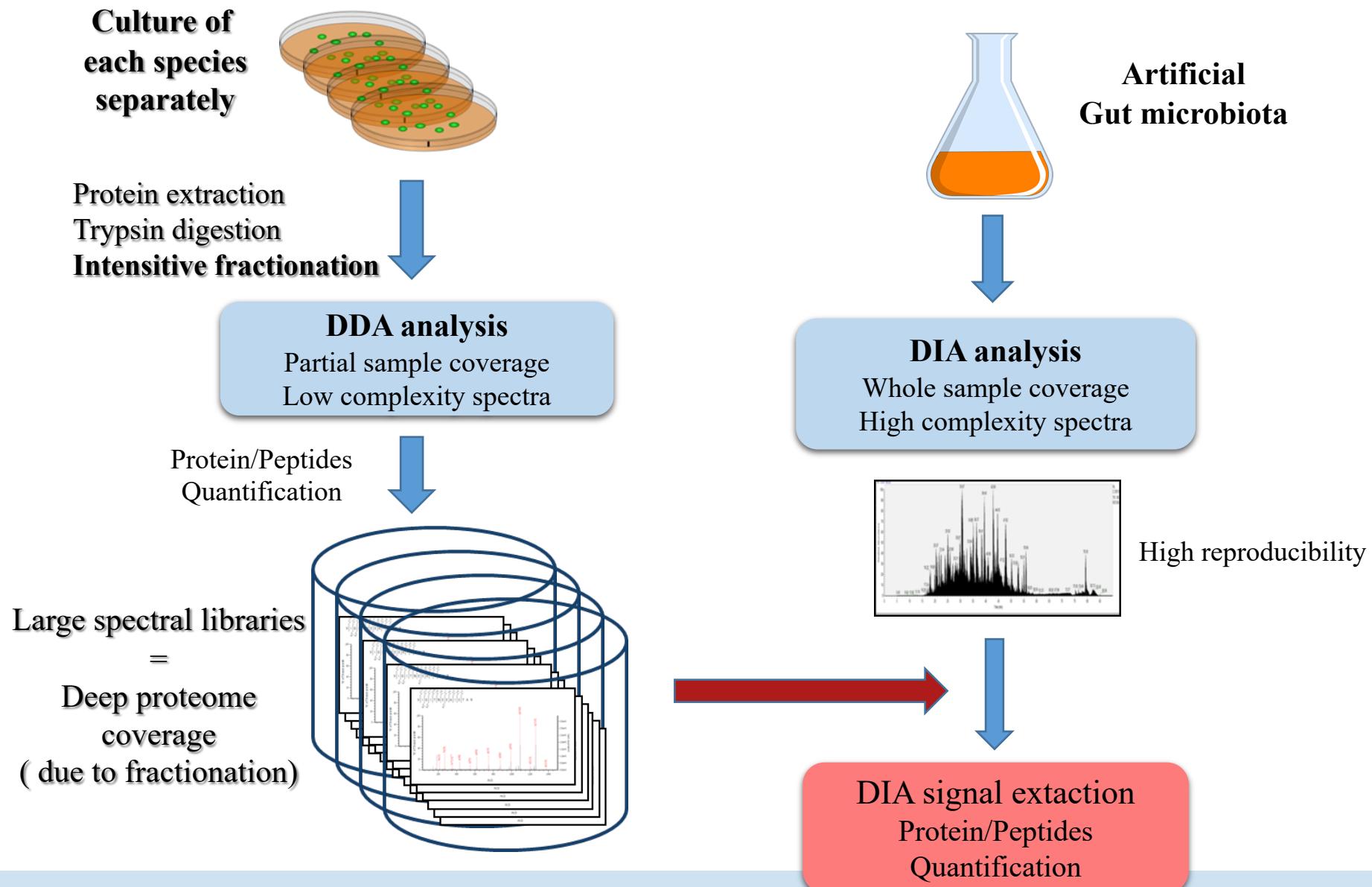


+ antibiotics  
(doses, time-course)

## Meta-Multi-Omics Analysis

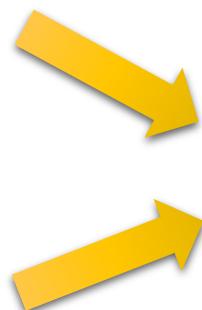
- Meta-Genomics
  - Meta-Transcriptomics
  - Meta-Proteomics
  - Meta-Metabolomics
- ]
- Data Integration

# Methods optimization for Deep coverage of Metaproteome

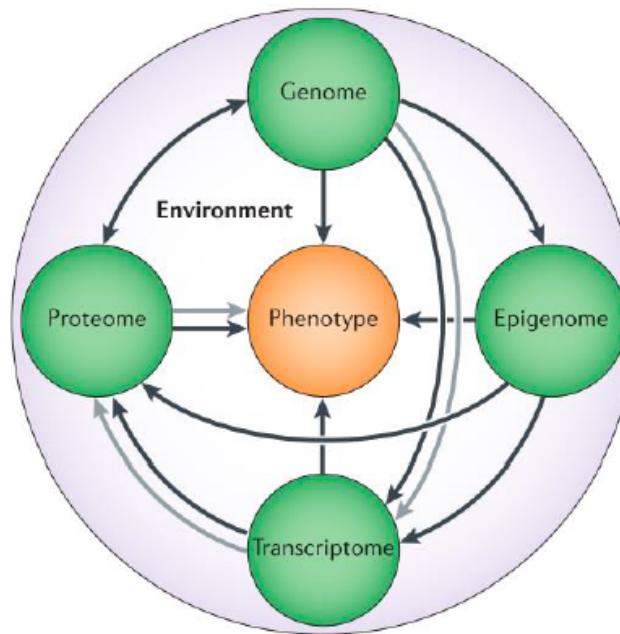


# multi-Omics integration

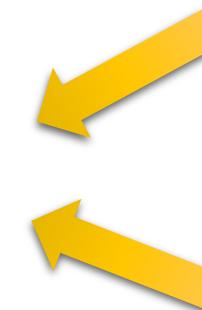
Global view of  
biological  
process



Complementary data



Complex and  
regulated at multiple  
levels



The more the  
better?

# Model for studying the effect of antibiotics on gut microbiota

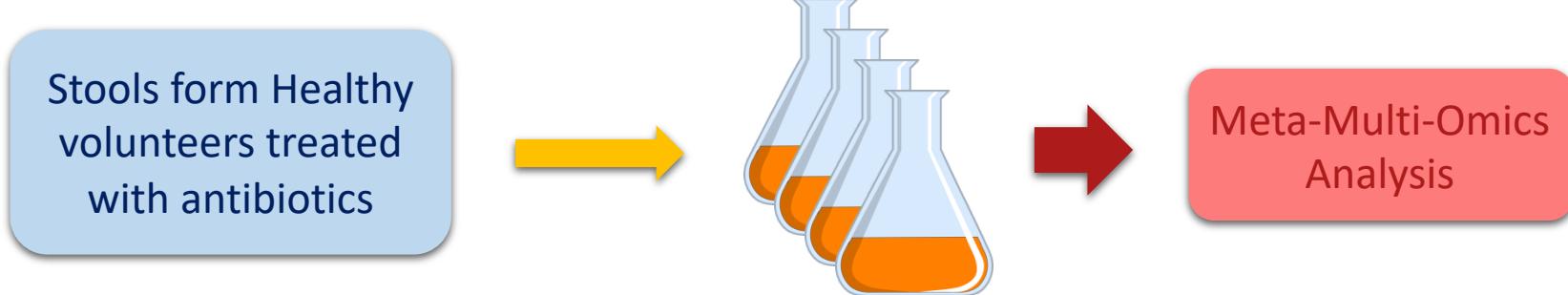
## Step1 : Model for studying the effect of antibiotics on gut microbiota



The artificial gut microbiota might be used for :

- Other studies (metabolomic diseases, effect of diets...)
- Improve reference databases for gut microbiota

## Step2: Application of developed multi-omics methods on « Culturomics »



## Step3: Direct Analysis from Patient specimens

Stools form Healthy volunteers treated with antibiotics



Meta-Multi-Omics Analysis

# Research Chair and Innovation

## L'Oréal in Digital Biology



CHAIRE DE RECHERCHE  
ET D'INNOVATION  
**L'ORÉAL**  
EN BIOLOGIE NUMÉRIQUE  
AFFILIÉ À  UNIVERSITÉ  
Laval

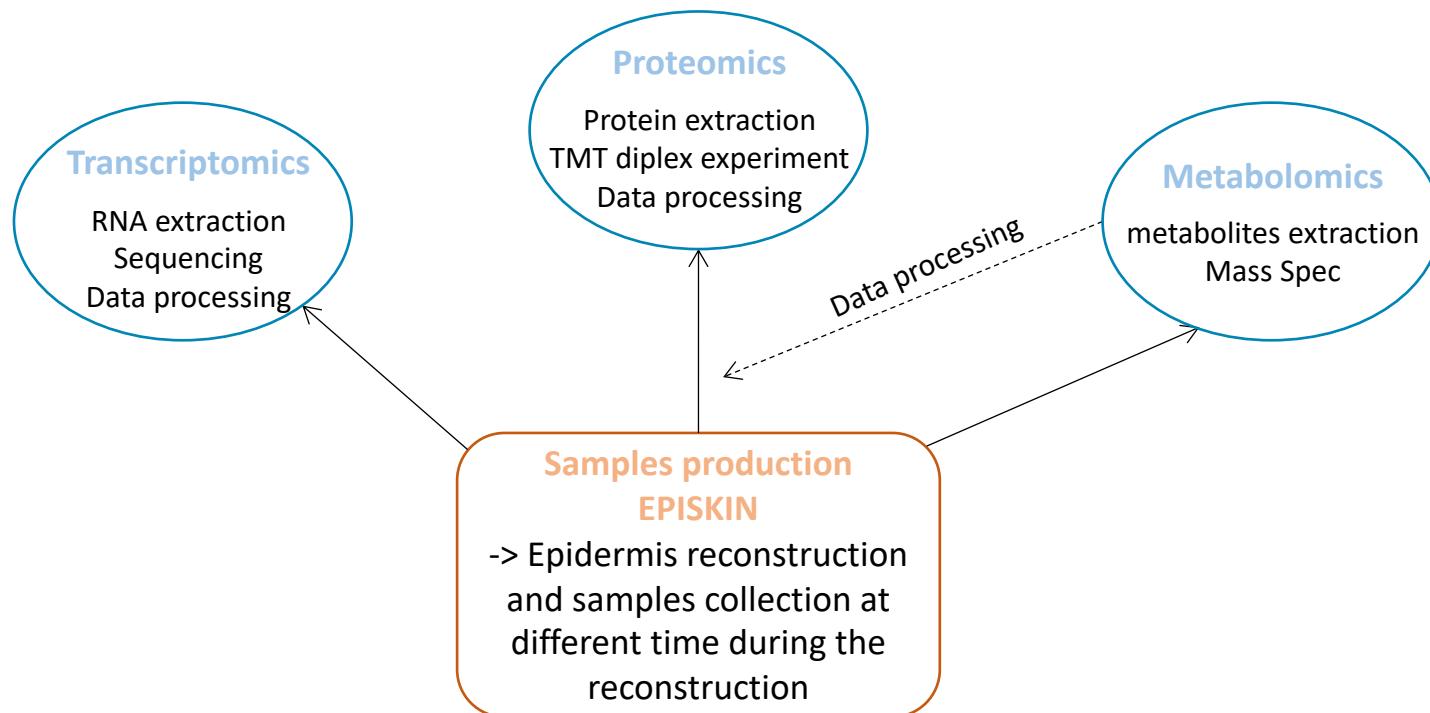
# Research Chair

- Epidermis reconstruction (Biological Knowledge)
  - Multi-OMICS analysis (Analytical strategies)
  - Co-Expression (Analytical strategies)
- Microbiome (Biological Knowledge – Analytical Strategies)
- Biomarker discovery by Machine learning
- Kibio.Science (Integrated Platform for big data)

# OBJECTIVES

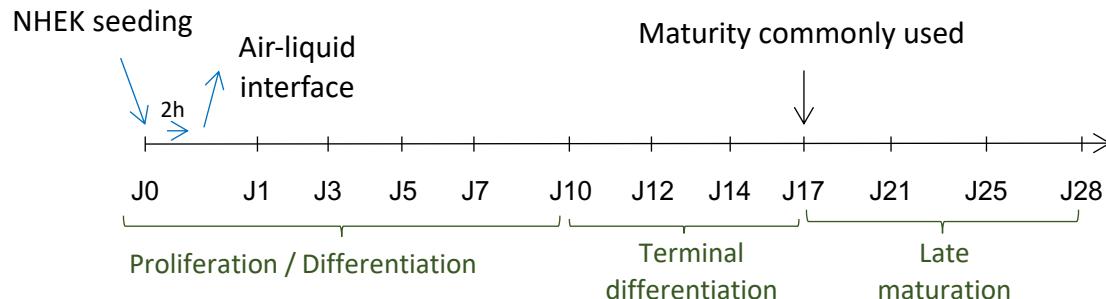
- Understand at the molecular level the epidermal differentiation during the reconstruction process, based on a multi-OMICS approach (transcriptomics, proteomics and metabolomics ), from J0 to J28
  - key biological players (genes, proteins, metabolites) involved in each step
  - key and transitional steps of the epidermal differentiation (pathways and biological funtions)
- Integrative biology: develop bioinformatics methods to correlate multi-OMICS signals

# Epidermis reconstruction (Collaboration with HTBD – L'Oréal)



## Samples preparation

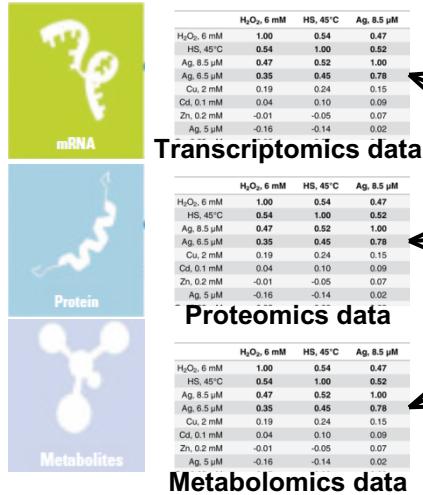
- **Reconstructed epidermis model => RHE (*Skinetic*)**
    - Cultured in a defined medium, without any medium changes during the reconstruction process
  - **Times chosen for the skin samples collection**



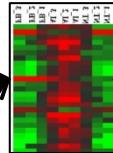
- **Replicates** -> Biological replicates rather than technical replicates => 4 donors
  - **Final design**

# Data analysis strategy

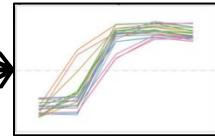
## Single Omics Analysis



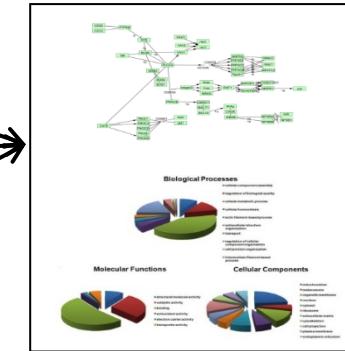
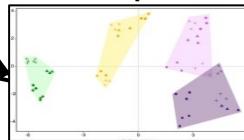
### Differential expression Analysis



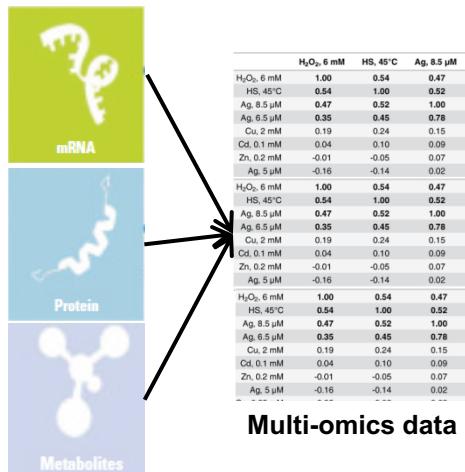
### Kinetic & co-expression Analysis



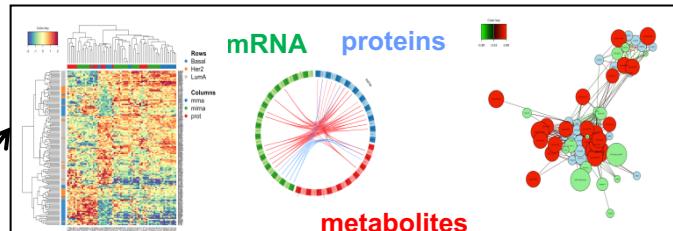
### Discriminant Analysis (PCA, PLS-DA, etc...)



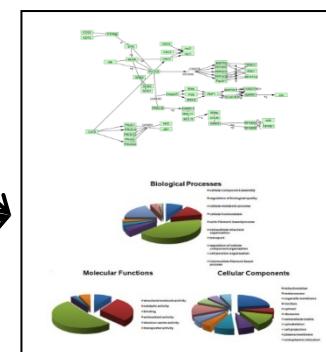
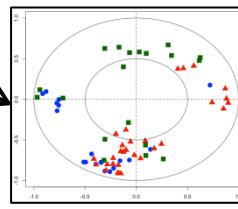
### Functional Analysis & Annotations (Pathway & Go enrichments)



## Multi-Omics Analysis



### Multi-Omics correlation analysis



### Functional Analysis & Annotations (Pathway & GO enrichments)

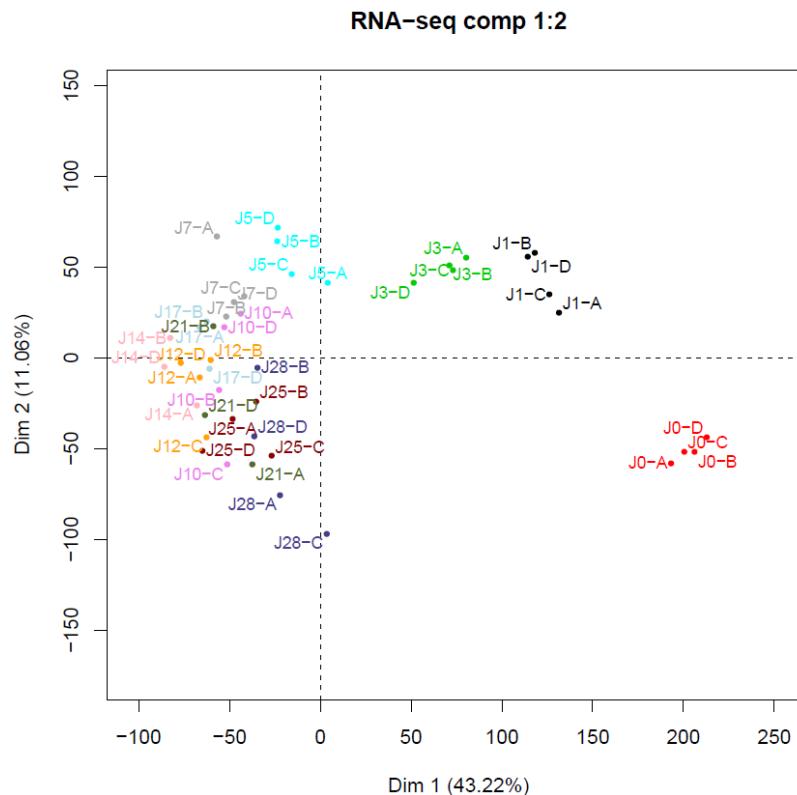
# Analytical Strategies

- Principal Component Analysis:  
→ Data quality and samples discrimination
- Kinetics Profiling  
→ Cluster of genes / proteins / metabolites with same kinetics profile
- Functional Analysis of Gene clusters
  - Identification of the main and successive biological functions involved during the epidermis reconstruction
- Transcriptomics Co-expression networks
- R Package MixOmics (exploratory phase)
  - Multi-Omics analysis
  - TimeOmics
- Machine learning

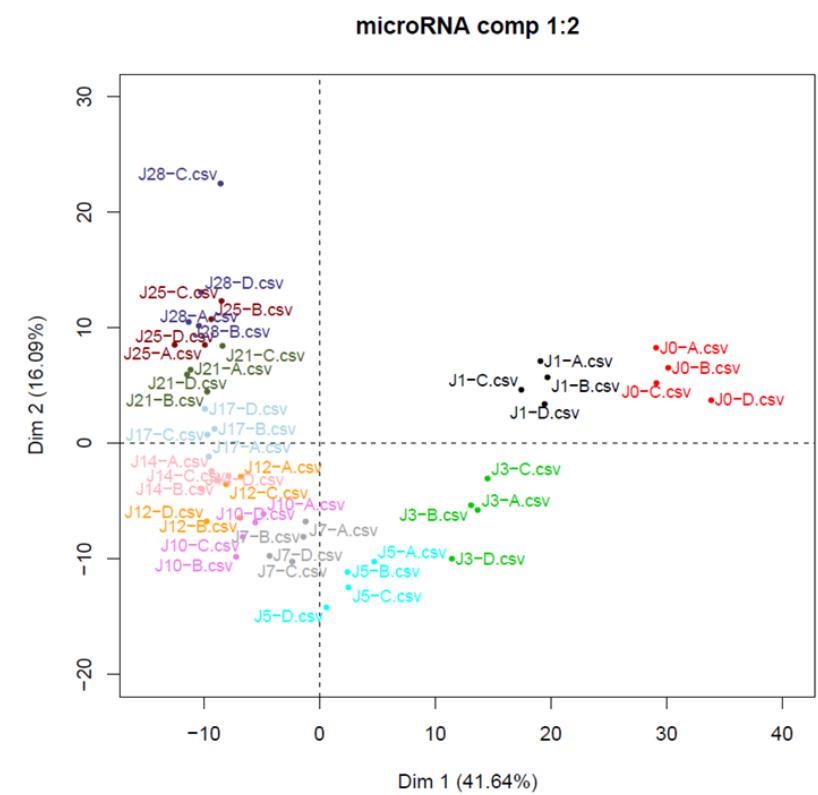
# Single omics analysis

miRNA-seq and total RNA-seq and the QC are achieved

43,038 transcripts (15,348 protein coding)

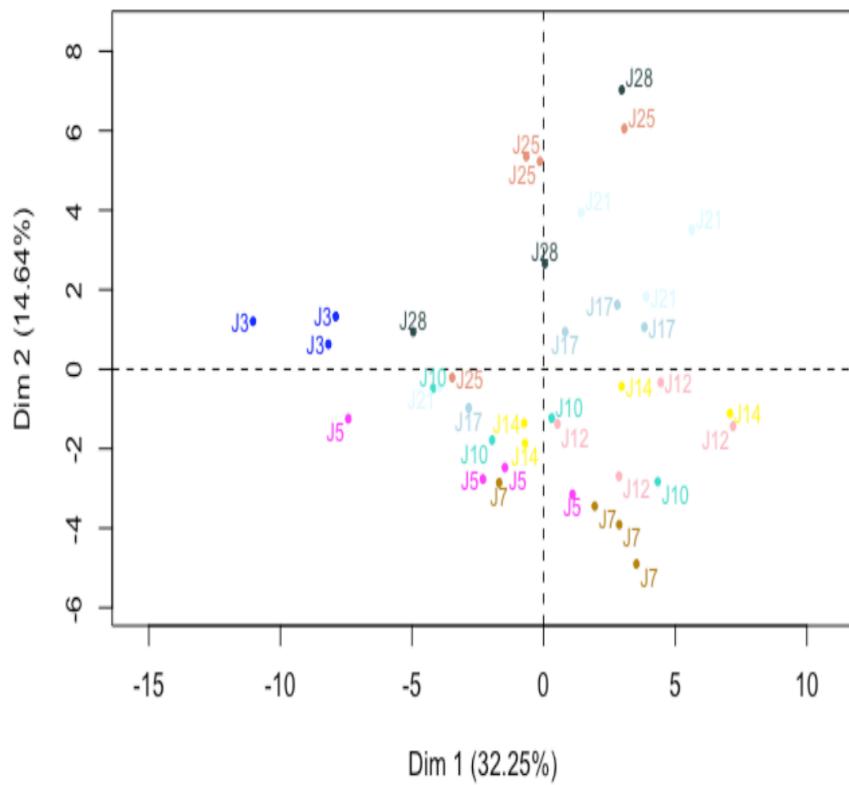


415 miRNA

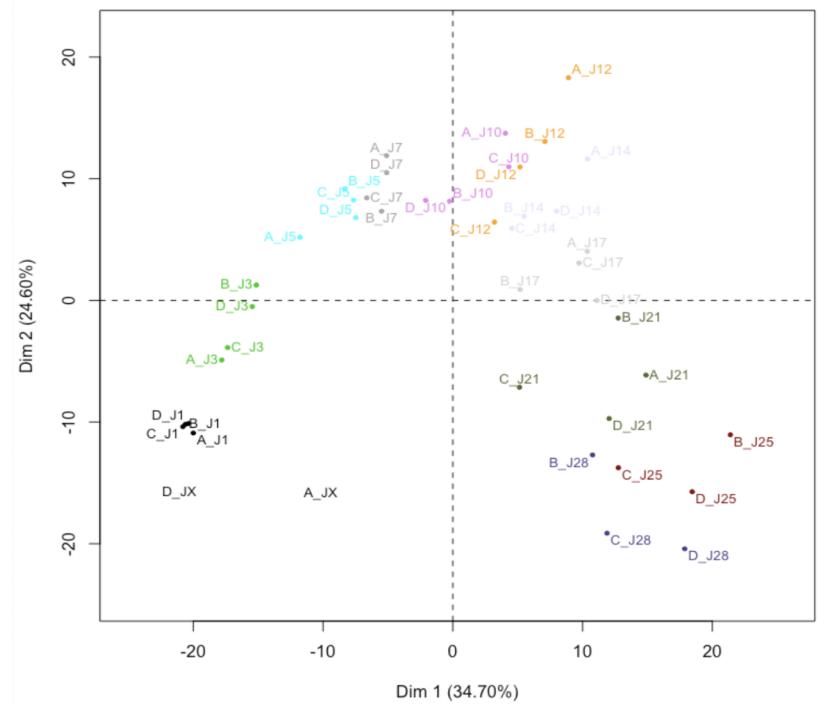


# Proteomics and metabolomics

Total of 1,820 quantifiable proteins (TMT)

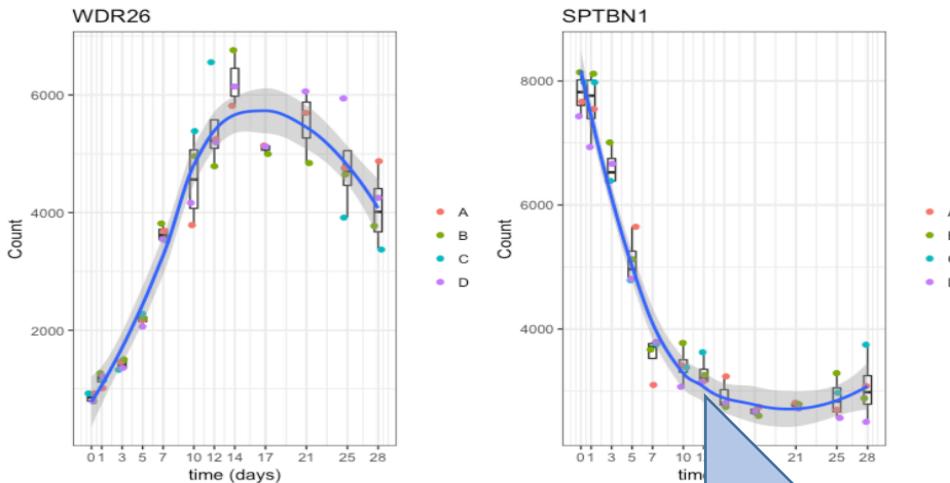


Total of 428 metabolites



# Time course analysis: transcriptomics

## Time course modelisation with splines (LMMS)

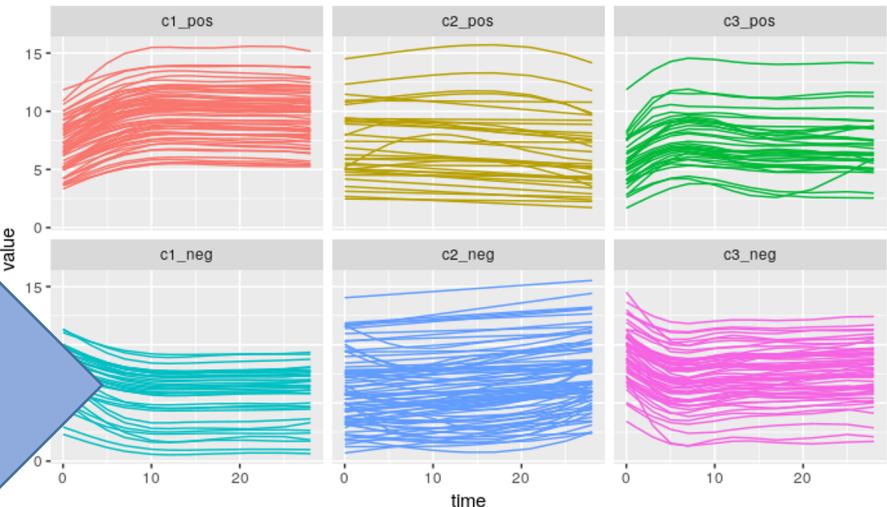


GO:0060429~epithelium development  
GO:0045595~regulation of cell differentiation

GO:0000278~mitotic cell cycle  
GO:0006461~protein complex assembly

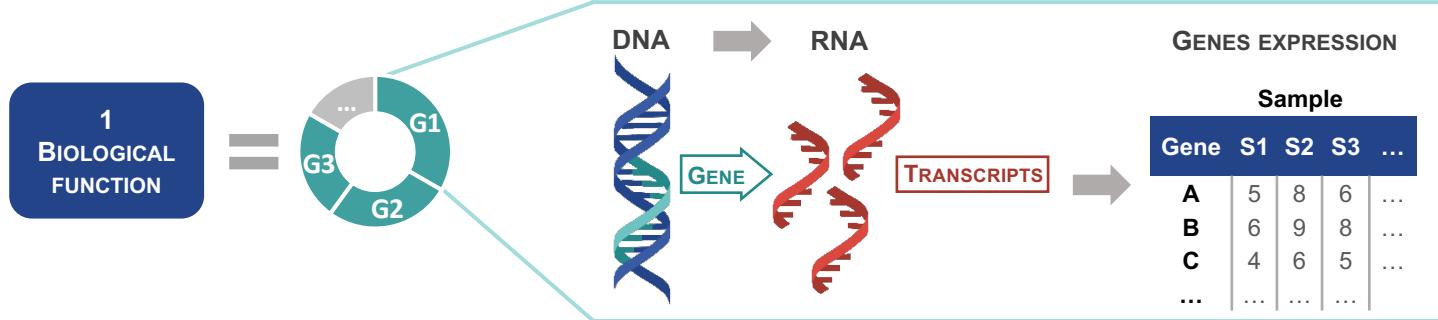
## Main clusters

sPCA clusters, keepX = c(100,100,100)

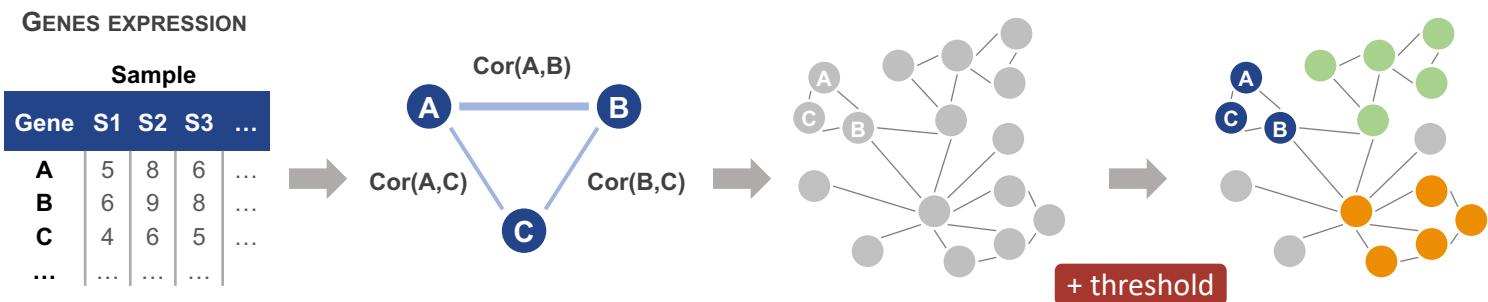


# Gene co-expression network

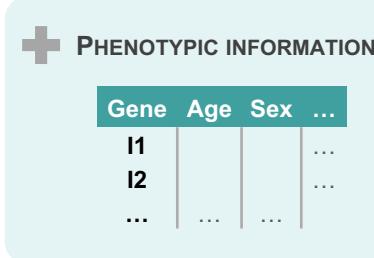
**BASE**



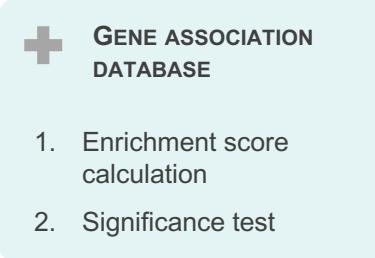
**CO-EXPRESSION  
METHOD**



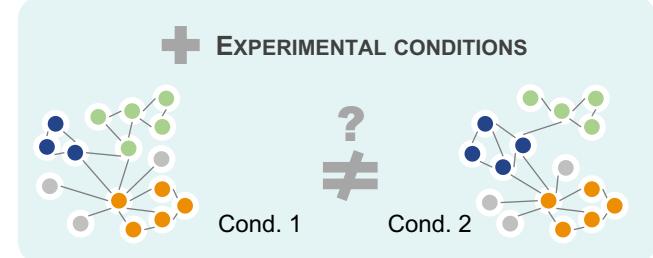
**EXPLOITATION**



**PHENOTYPIC ASSOCIATION**



**FUNCTIONAL ENRICHMENT**



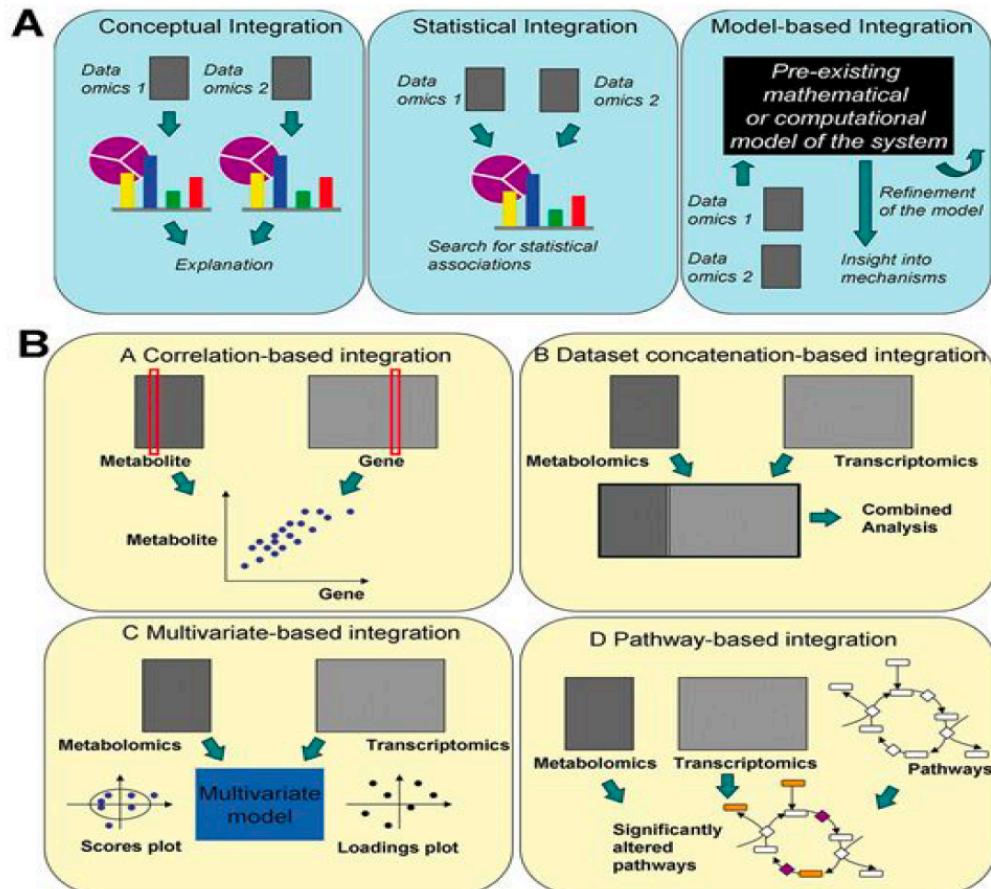
**TOPOLOGICAL DYNAMIC**

# Multi-Omics analysis

# Simple multi-Omics Kinetics Profiling

- Multi-Omics Kinetics profiling
- Does genes and corresponding protein have same kinetics profile ?  
→ Do we enrich same biological function ?
- mRNA-miRNA target predictions based on experimental experiments

# Advanced Multi-Omics approaches Overview



- Static integration
- Time course integration

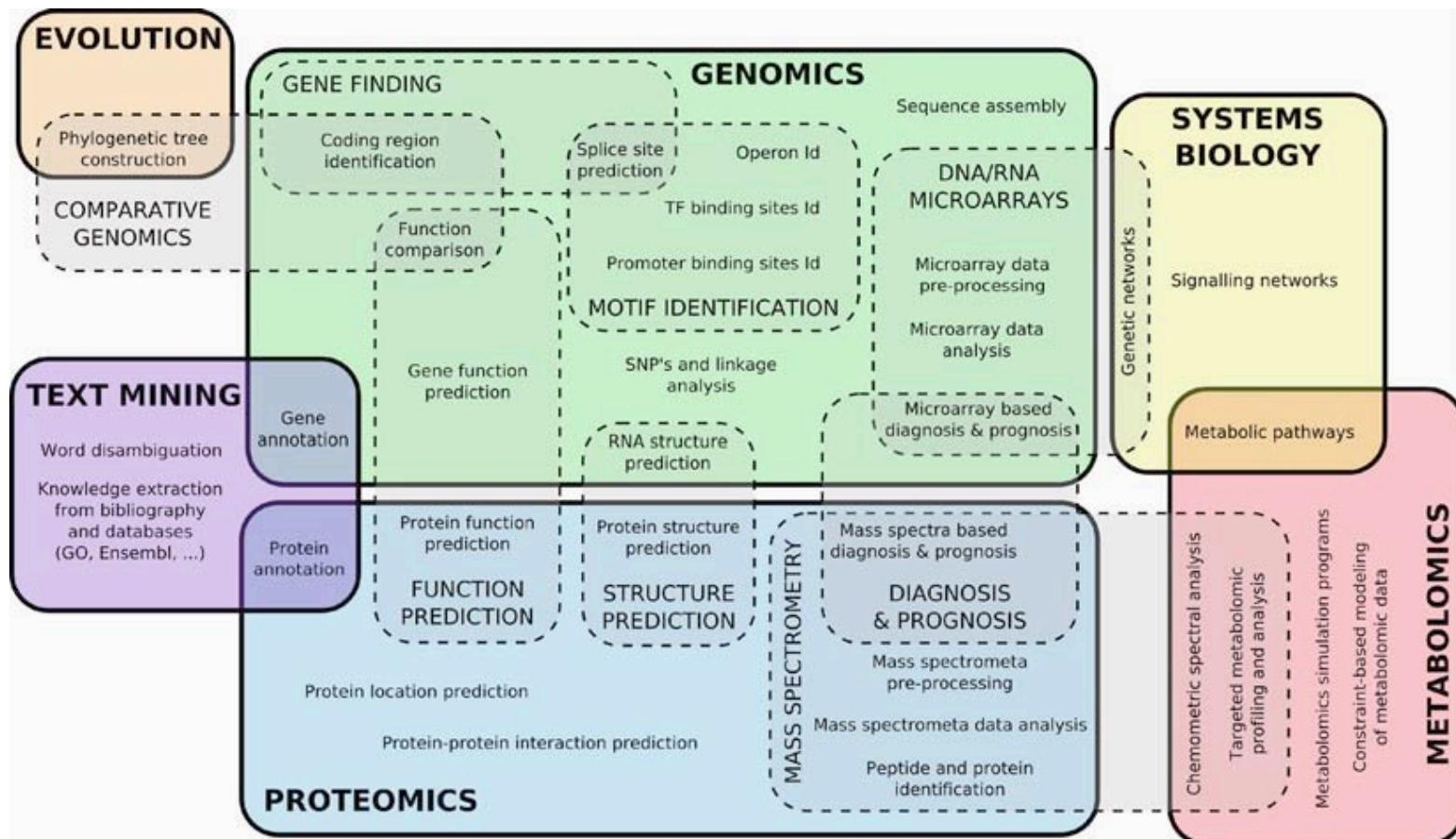
Cavill, R., Jennen, D., Kleinjans, J., & Briedé, J. J. (2016). Transcriptomic and metabolomic data integration. *Briefings in bioinformatics*, 17(5), 891-901.

# Next steps & Perspectives

- Continue single omics analysis
  - Kinetics profiling
  - Main biological functions
  - Discriminate molecular entities for each maturation day
- Multi-Omics : new strategies to identify innovative biomarkers
  - Improve the robustness of molecular signatures
  - Integrate complexity of biological event
    - Complex regulation between multiple biological entities (mRNA, miRNA, proteins, metabolites)
    - Time lags
- Improve our knowledge of the epidermis formation of in vitro model by deciphering the successive events from transcripts to the metabolites
  - Key events
  - Key molecular entities
  - identify more robust multiple targets for developing new cosmetic concepts

Identify biomarker signatures in  
OMICs data using machine learning  
approaches

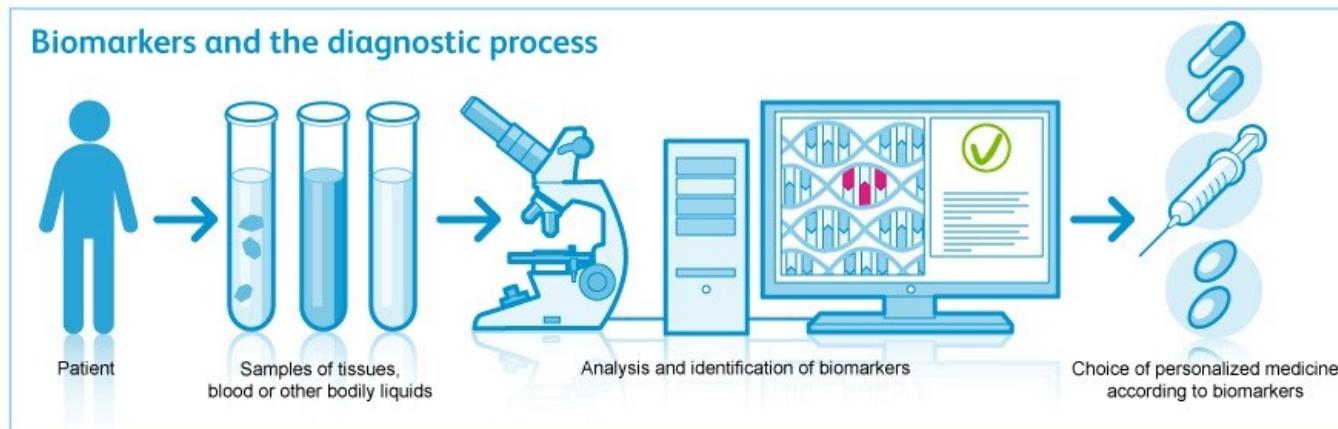
# Machine learning applications in biology



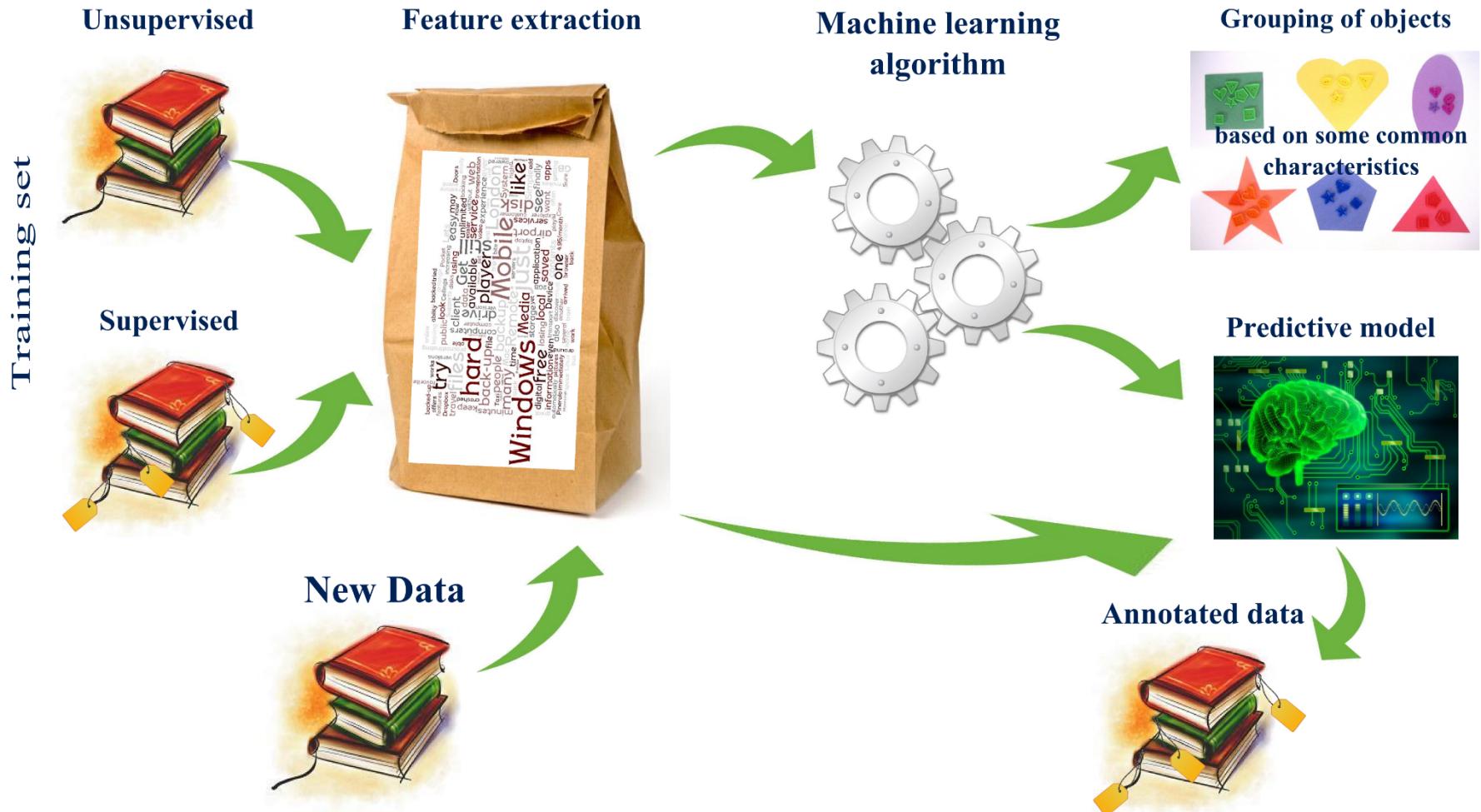
Inza et al, 2010

# Biomarker signature discovery

- Biomarker = biological marker
- Substance, structure, or process that can be measured in the body or its products and influence or predict the incidence of outcome or disease
- Indicator of medical state observed from the patient

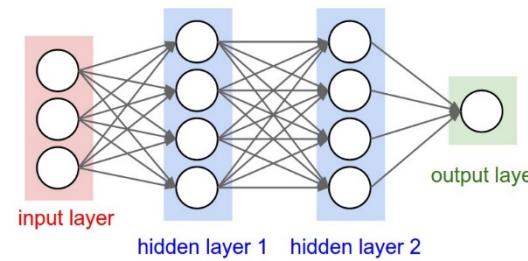
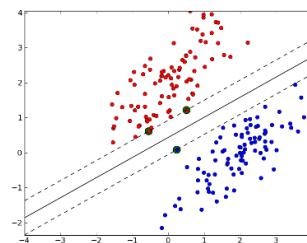


# Machine learning workflow



# Challenges and Big questions

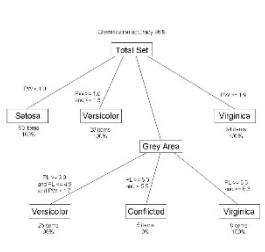
- OMICs: Tens of thousands of features vs 20~1000 samples
- What feature ranking/selection algorithms to use?
  - Univariate: Information gain, Gain ratio, RELIEFF...
  - Multivariate: Correlation based, Markov blanket filter...
  - What is the number of best features to include in the model?
- What classification/regression algorithms to choose?
  - Decision trees, Random Forest, SVM, NeuralNets, NaiveBayes...
- What hyper-parameters to set?



$$P(c | X) = P(x_1 | c) \times P(x_2 | c) \times \dots \times P(x_n | c) \times P(c)$$

$$P(c | x) = \frac{P(x | c)P(c)}{P(x)}$$

Likelihood                      Class Prior Probability  
Posterior Probability           Predictor Prior Probability



# Proposed approach

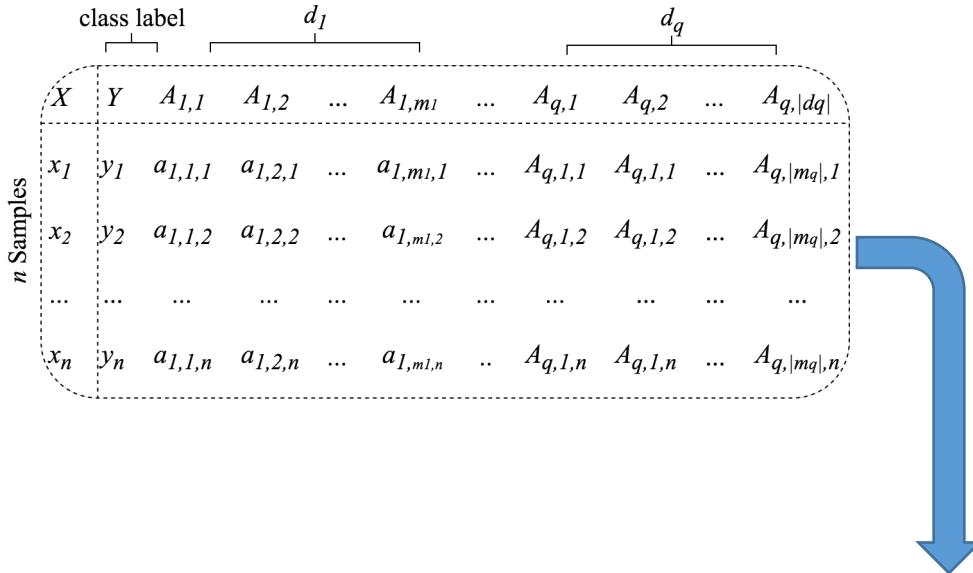
**Objective: Help in the choice of the classifier and the best subset of features**

- Tests all algorithms available in the weka library
- Multi-threaded supervised forward stepwise feature selection

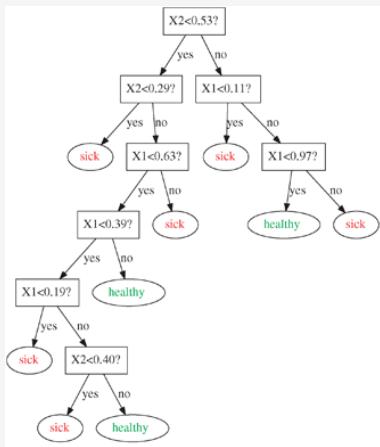
➤ BioDiscML  
(Biomarker Discovery by Machine Learning)

# Input / output

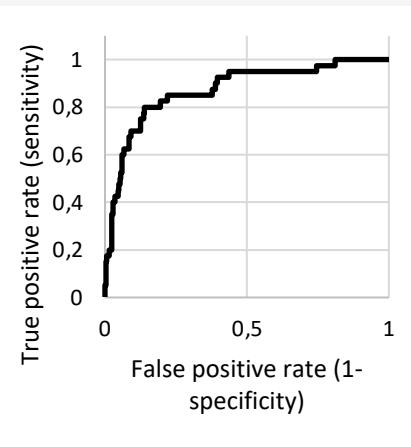
## Input



## Output



## Models



## Performances

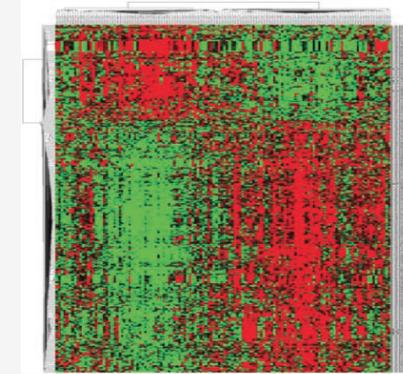
## Samples (e.g. patients)

## Features

- Gene expression
- Non coding RNA expression
- Protein peptides quantification
- Methylation
- Polymorphisms
- Metabolites quantification
- Clinicopathologic data

## Class

- Outcome(s) (e.g. disease vs normal)
- Category (e.g. bacteria species, cancer stage)
- Level (e.g. product concentration)
- Time (e.g. cancer recurrence)



## Signatures

# Supported ML classifiers

About 85 classifiers, various compatibilities

Category	Classifier	Supported feature types			Supported class types			Ref.
		Nom	Num	MV	Nom	Bin	Num	
Bayes	Averaged N Dependence Estimators (A1DE, A2DE)	X	X	X	X	X		(77)
	Bayes Network	X	X	X	X	X		(78)
	Complement Naive Bayes		X	X	X	X		(79)
	Hidden Markov Model	X	X		X	X		(78, 80)
	Hidden Naive Bayes	X			X	X		(81)
	Naive Bayes	X	X	X	X	X		(82)
	Multinomial Naive Bayes		X		X	X		(83)
	Simple Naive Bayes	X	X	X	X	X		(84)
Functions	ElasticNet		X				X	(85)
	Fisher's Linear Discriminant (FLDA)		X				X	(86)
	Gaussian Processes	X	X	X			X	(87)
	Isotonic Regression		X				X	(88)
	Kernel Logistic Regression	X	X	X			X	(89)
	Latent Dirichlet allocation (LDA)		X		X	X		(90)
	Least Median Squared linear regression	X	X	X			X	(91)
	Liblinear	X	X	X	X	X		(92)
	Support Vector Machine LibSVM	X	X	X	X	X		(93)
	Linear Regression	X	X	X			X	(94)
	Multinomial Logistic Regression	X	X	X	X	X		(95)
	Multilayer Perceptron (multiple implementations) (MLP)	X	X	X	X	X	X	(96)
	Quadratic Discriminant Analysis (QDA)		X		X	X		(97)
	Radial Basis Function (RBF), incl. Network and regressor	X	X	X	X	X	X	(98)
	Stochastic Gradient Descent (SGD)	X	X	X		X		(99, 100)
	Simple Linear Regression		X	X			X	(101)
Trees	Simple Linear Logistic regression	X	X	X	X	X		(102)
	Sequential Minimal Optimization (SMO)	X	X	X	X	X		(103–105)
	Sequential Minimal Optimization regression	X	X	X			X	(106, 107)
	Stochastic Primal Estimated sub-GrAdient SOlver for SVM (Pegasos)	X	X	X		X		(108)
	Voted Perceptron	X	X	X		X		(109)
	Winnow		X			X		(110)
	Analogical Modeling	X			X	X		(111)
	K nearest neighbours classifier with various distance measures	X	X	X	X	X		(112, 113)
	CHIRP						X	(119)
	Fuzzy Lattice Reasoning (FLR)						X	(120)
Misc.	HyperPipes						X	(40)
	Ordinal Stochastic Dominance Learner (OSDL)						X	(121, 122)
	Voting Feature Intervals (VFI)						X	(123)
	ConjunctiveRule						X	(40, 124)
	Decision Table/Naive Bayes hybrid (DTNB)						X	(125)
	Fuzzy Unordered Rule Induction Algorithm (FURIA)						X	(126)
	Repeated Incremental Pruning to Produce Error Reduction (J Rip)						X	(127)
	Lazy Associative Classifier (LAC)						X	(128)
	M5 Rules						X	(129–131)
	MODLEM						X	(132)
	Multi Objective Evolutionary Fuzzy						X	(133)
	Nearest-neighbor-like generalized (NNge)						X	(134)
	Ordinal Learning Method (OLM)						X	(135)
	One Rule (OneR)						X	(135, 136)
	PART						X	(137)
	PRISM						X	(133, 138)
	Ripple-DOWN Rule (Ridor)						X	(139)
Random Forest	Rough Set						X	(112)
	Zero Rule (ZeroR)						X	(39)
	Alternating Decision Tree (ADTree)						X	(112, 140)
	Alternating Model Tree						X	(141)
	Best-First Decision Tree (BFTree)						X	(142)
	Credal Decision Tree (CDT)						X	(143)
	Decision Stump						X	(144)
	ExtraTree						X	(145)
	Functional Trees (FT)						X	(146)
	Hoeffding Tree						X	(147)
Ensemble	ID3						X	(146, 148)
	C4.5 Tree (J48)						X	(149)
	C4.5 Consolidated Tree (J48 Consolidated)						X	(150, 151)
	LogitBoost Alternating Decision Tree (LADTree)						X	(152)
	Logistic Model Trees (LMT)						X	(102, 146)
	M5 Model Trees						X	(130, 131)
	Naive Bayes Tree (NBTree)						X	(153)
	Random Forest						X	(146)

# Use case : Prostate cancer

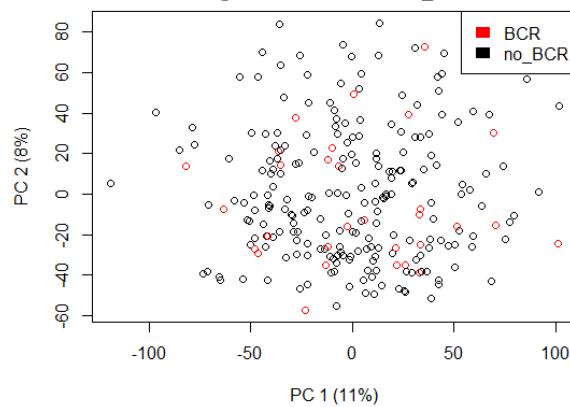
**Goal:** Predict the biochemical recurrence in prostate cancer

**Input:** TCGA, 329 patients. Gene expression (RNAseq, 19k genes)

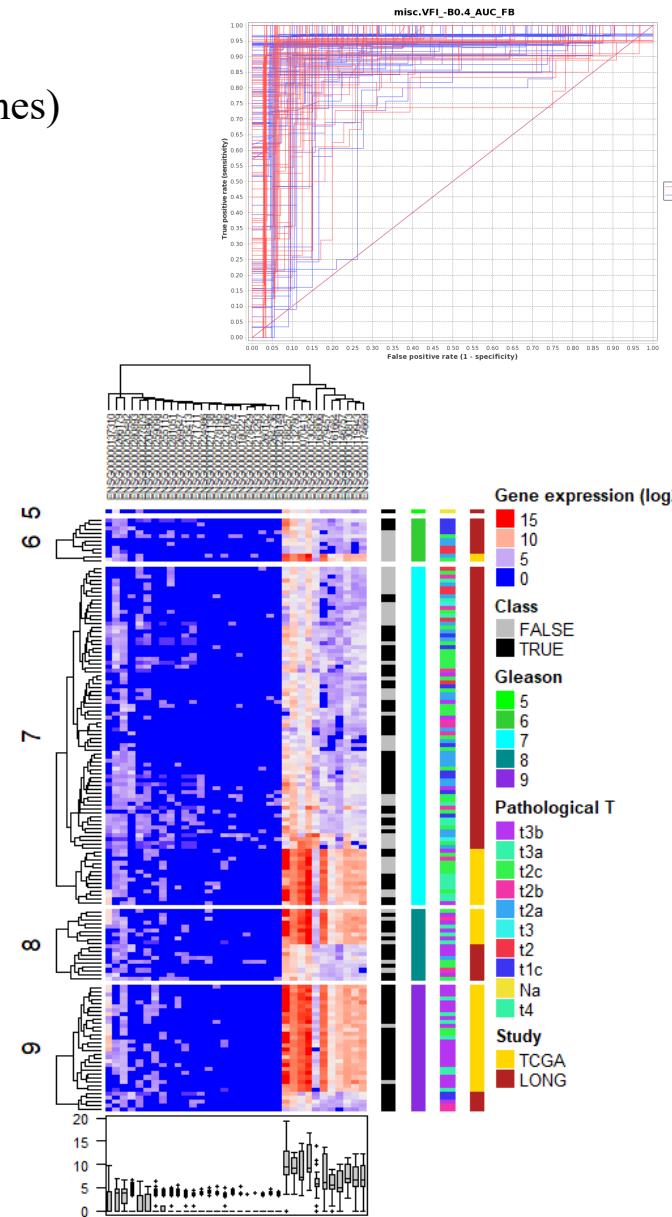
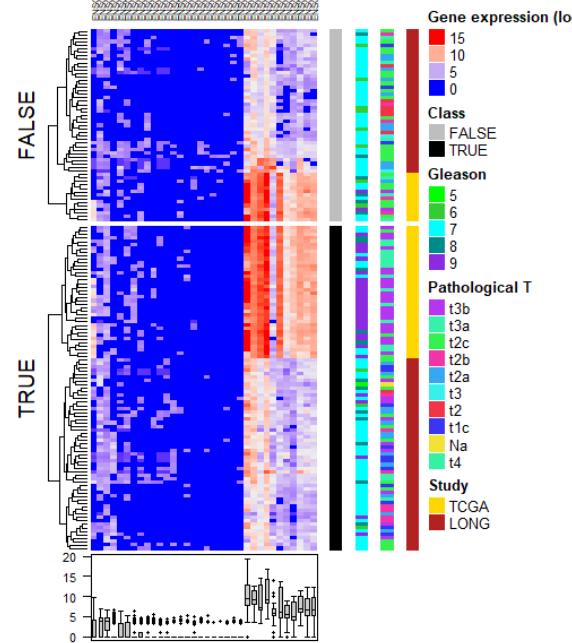
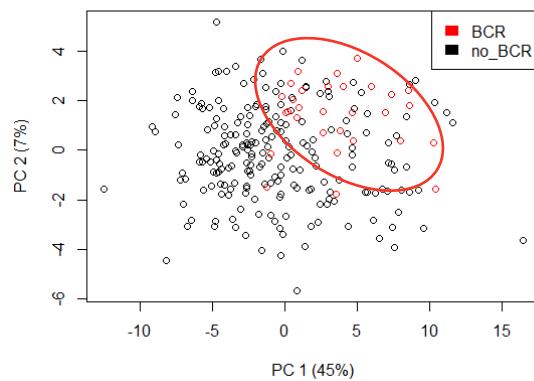
**Model:** voting features intervals (bias 0.4)

**AUC:** 0.8 (bootstrapping)

All genes in input

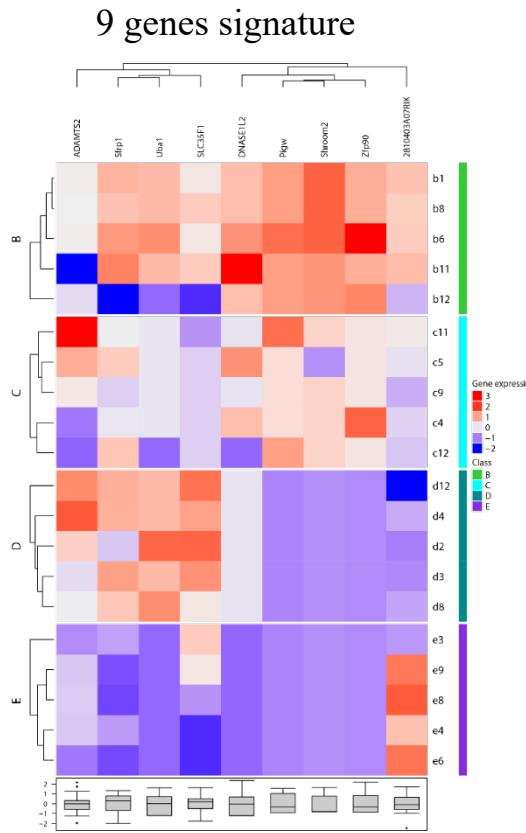


Signature (31 genes)

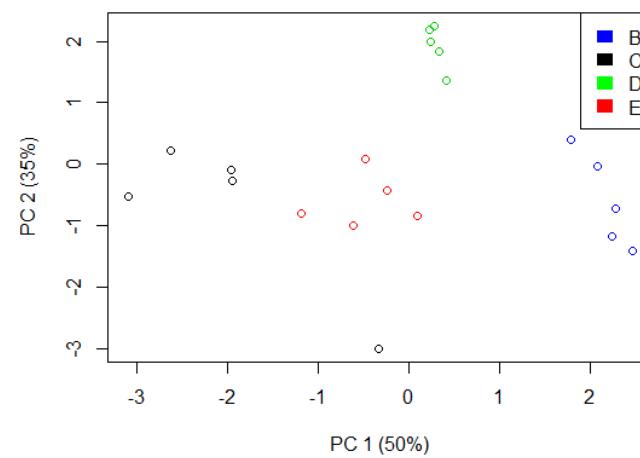
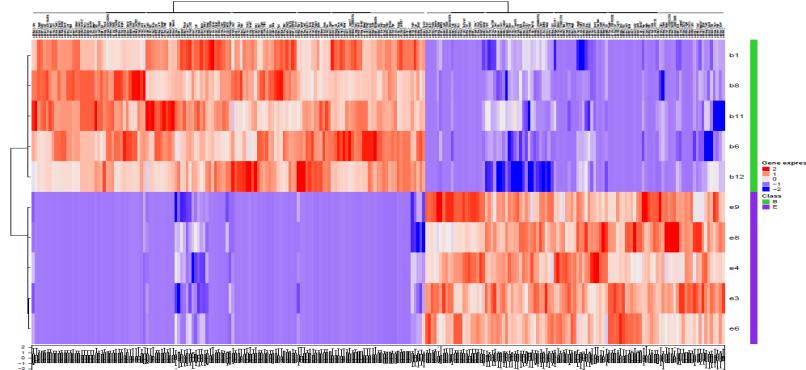


# Use case: Mice fed by different fatty acids

- **Study:** 4 groups of mice have been fed different fatty acids
- **Goal:** Find a gene signature that discriminate the mice groups
- **Input:** Gene expression (RNAseq, 22k features)
- **Model found:** K-nearest neighbours classifier
- **AUC:** 0.93



Correlated genes by spearman correlation (classes B vs E) and 10CV infogain merit >0.9:



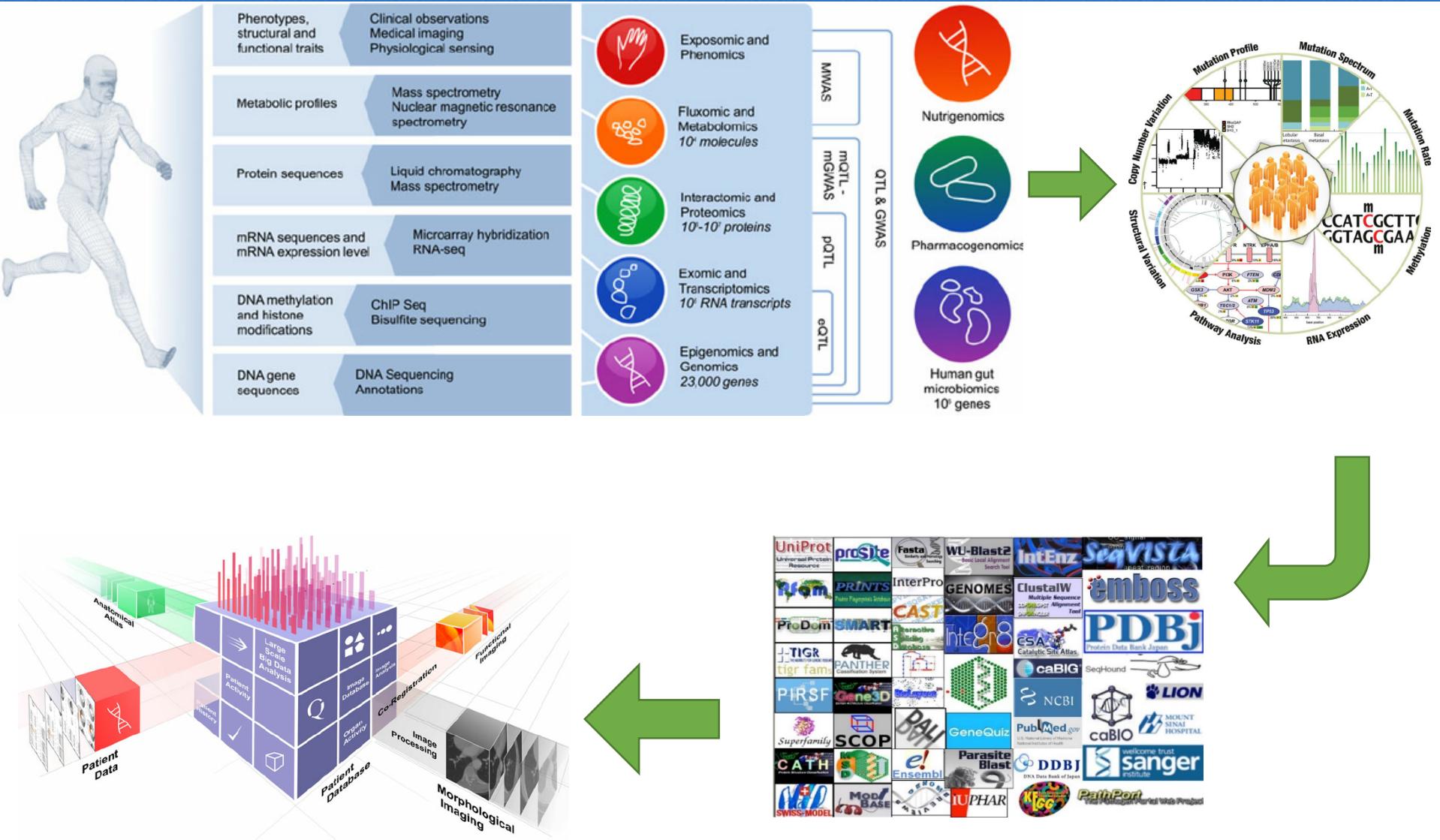
# Conclusions

- BioDiscML is an automated approach to search for best ML algorithms which fit input data
- Help use to decide what model to use
- Identifies the best subset of features for prediction
- Report the correlated genes, useful to understand

# Kibio.science



# From patient to massive data



Andreu-Perez et al. 2015

# Evolution for databases

- Evolution in molecular biology
  - From 2 databases in 1980<sup>[1]</sup> to 1737 databases in 2017<sup>[2]</sup>
- Databases of biological databases
  - MetaBase<sup>[1]</sup>, NAR<sup>[2]</sup>, NIST<sup>[3]</sup>, OReFiL<sup>[4]</sup>, DBD<sup>[5]</sup>, HSLS<sup>[6]</sup>, and more

[1] D. J. Rigden et al., Nucleic Acids Research, 2012

[2] D. J. Rigden and X. M. Fernández, Nucleic Acids Research, 2018

[3] D. M. Blakeslee et al., NIST Data Gateway, 2016

[4] Y. Yamamoto and T. Takagi, BMC Bioinformatics, 2007

[5] M. V. N. Setty and H. Rao, Jour. of Biochemical Technology, 2009

[6] A. Chattopadhyay et al., J Med Libr Assoc, 2006

# Challenges for a researcher

- Many databases & many web services

- Need to...

- Jump from one website to another to get dataset
- Use external tools to create relations
- Use external tools to analyse
- Visualization



1/ Hard  
2/ Time-consuming  
3/ Repeat

# Kibio project

- Huge need to connect biological knowledge together
- Link data easily
- New search engine
- Online portal



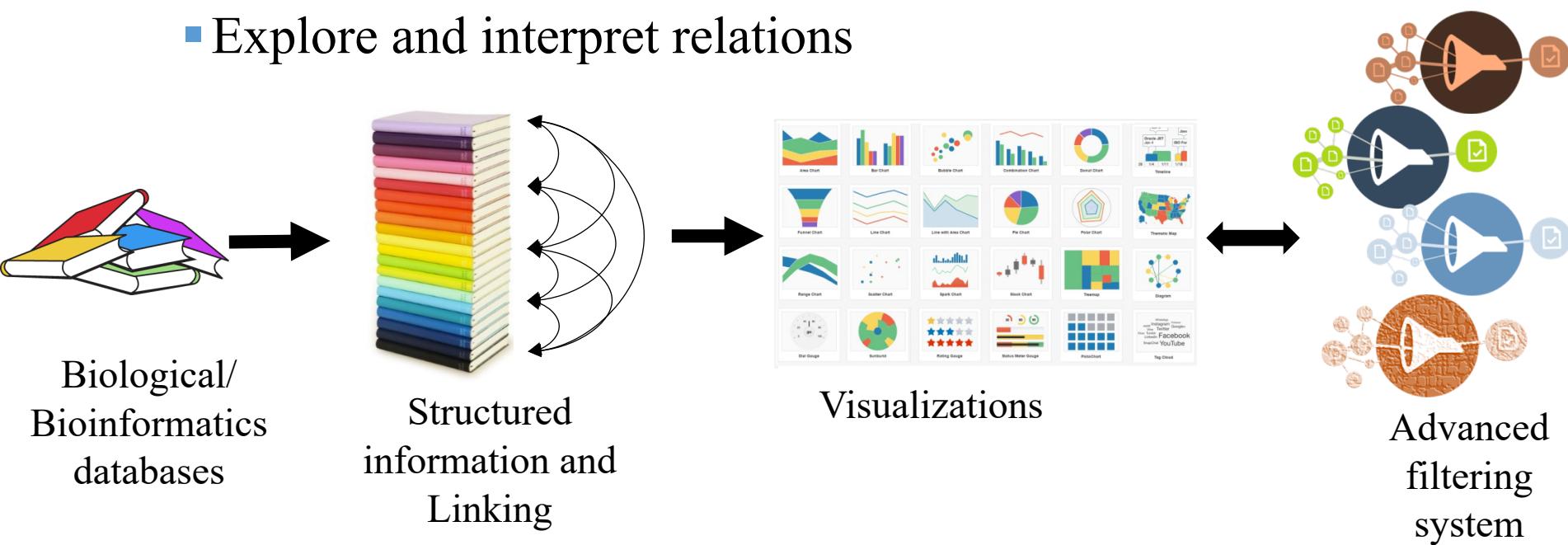
# Challenges to manage the data

- How to efficiently store?
    - Volumes,
    - Data structures,
    - Links,
    - Etc
  - How to quickly explore and search?
    - Splitting,
    - Algorithms,
    - Requests,
    - Etc
- 
- Interconnected
-  elasticsearch
-  SIREN  
DATA INTELLIGENCE

# Objectives

## Goal: Create a platform to

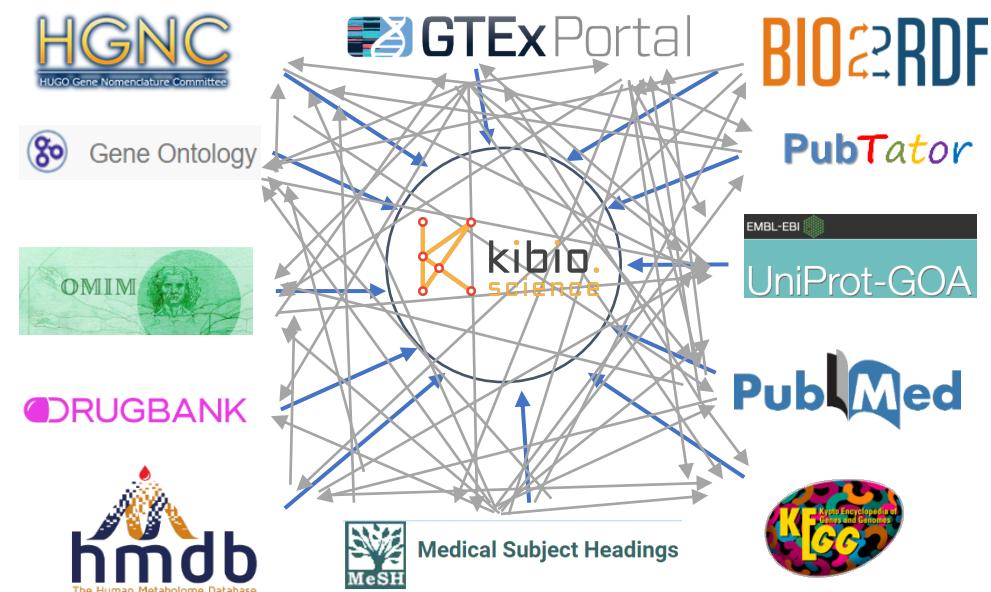
- Store and join together most known bioinformatics databases at one place
- Provide various and flexible plugins to visualize life science data
- Perform instant searches and interact dynamically with the visualizations
- Explore and interpret relations



# Database integration and linking



With relations

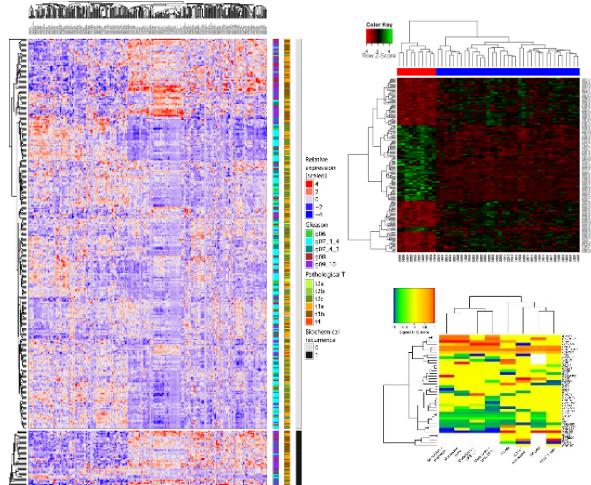


# Visualizations

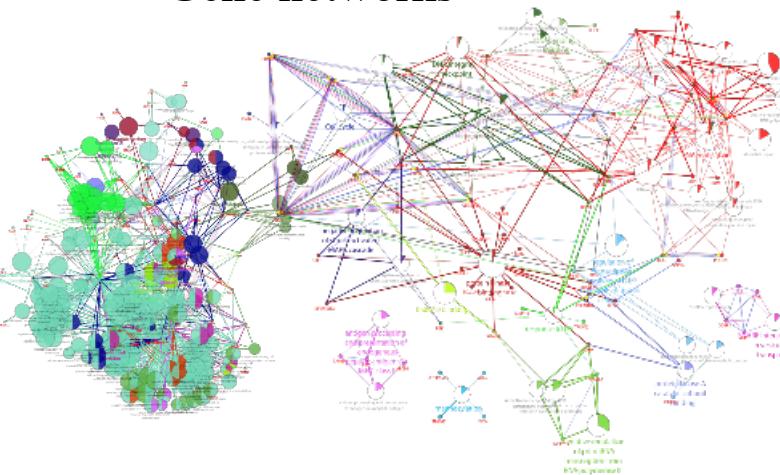


# More visualizations

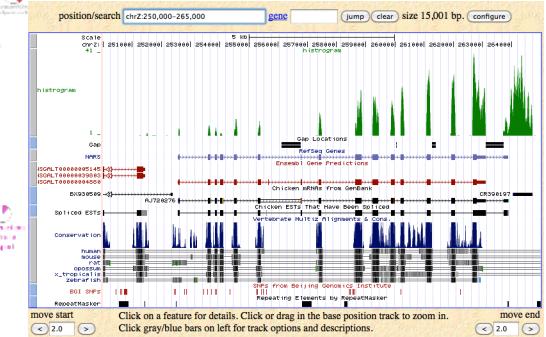
## Expression heatmap



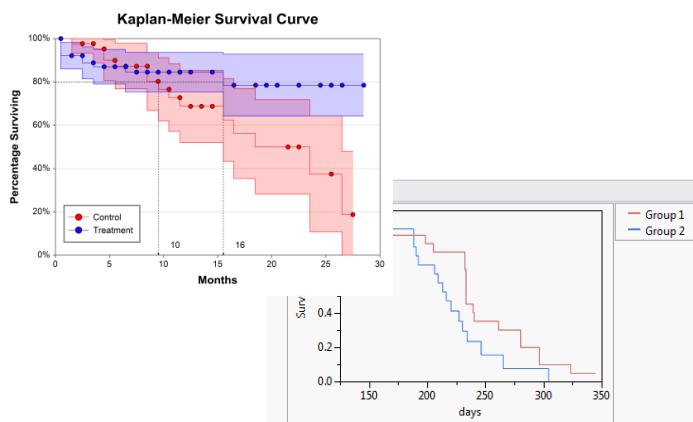
## Gene networks



## Genome browser

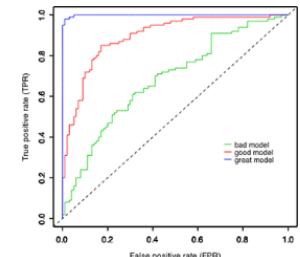


## Survival analysis



## Biostatistics

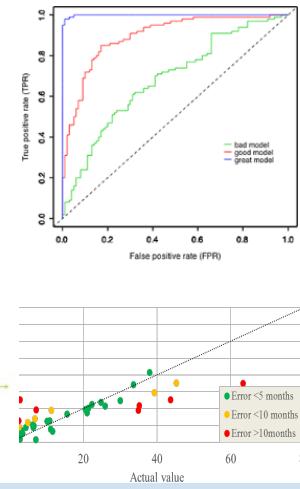
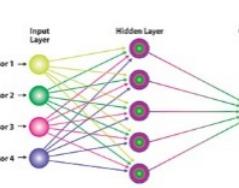
Model Variable	All Cases Mortality			CV Mortality		
	Conditional Probability of death #1	Probability of death #2	Conditional Accuracy (% index)	Conditional Probability of death #1	Probability of death #2	Conditional Accuracy (% index)
DTS	-0.001	0.001	53.9	0.005	-0.001	51.8
HS, age, history of congestive	0.001	0.001	53.9	0.005	0.001	51.8
HS, prior coronary disease	0.001	0.001	53.9	0.005	0.001	51.8
HS, age, history of congestive, HS, prior coronary disease, HS increase	0.001	176.4	0.746	0.002	174.2	0.833
HS, age, history of congestive, HS, prior coronary disease, HS increase	-0.001	175.6	0.741	-0.001	153.4	0.824
HS, age, history of congestive, HS, prior coronary disease, HS increase, HS, prior coronary disease, HS increase	0.003	177.6	0.745	0.003	138.0	0.824



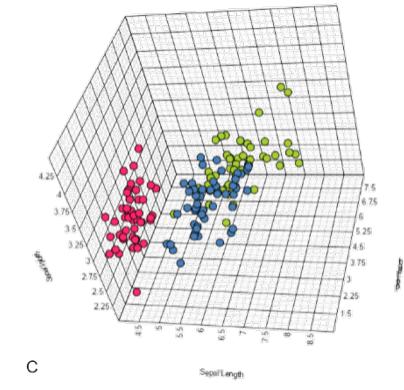
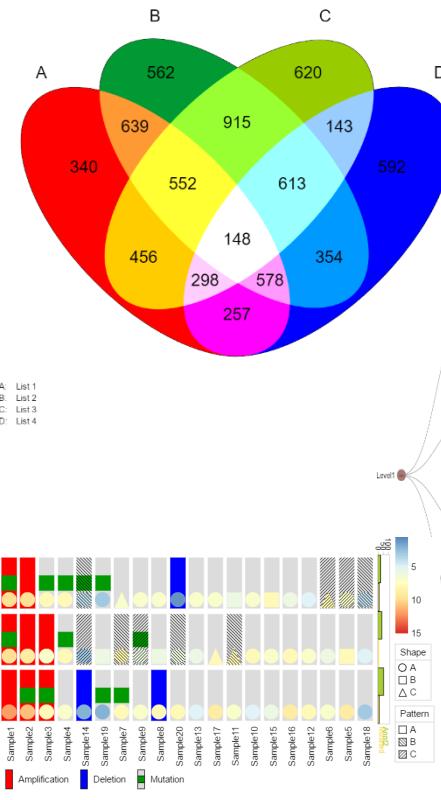
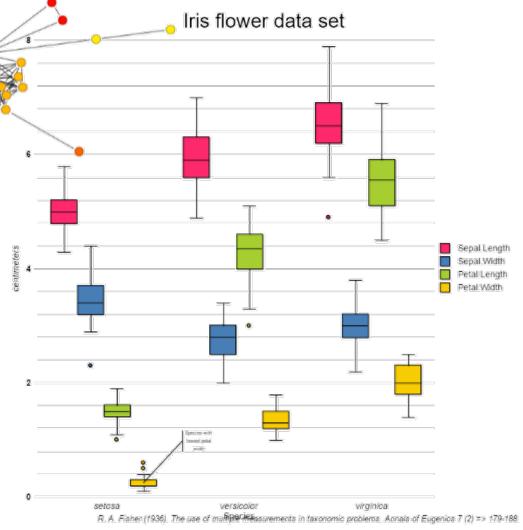
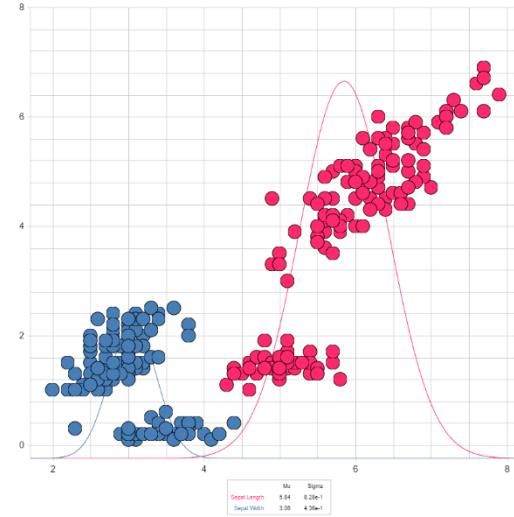
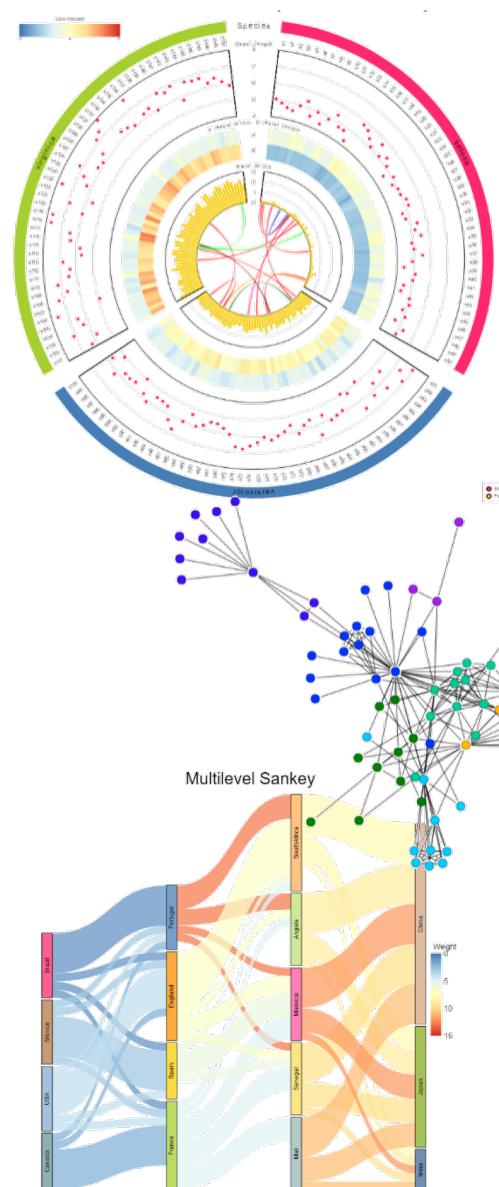
## Machine learning

Variable	GENDER		RSESSES				
	Low exp.	High exp.	Low exp.	High exp.			
All	487 81%	12 17%	358 77%	111 27%			
Age (years)	43.45	33.27	26 22%	6.75	47 14%	29 8%	6.02
Female (F)	37.45	15.9	32 35	4.75	71 15%	122 9%	4.75
HS (yes)	65.78	119	24.21	30%	44 9%	86 7%	30%
0-17	332	72.5	71 87%	0.829	161 33%	262 20%	6.716
18-49	46	9.5	10 12%	1.78	38 8%	38 3%	1.78
50-64	206	42.5	21 24%	1.78	38 8%	38 3%	1.78
65+	59	8.5	5 7%	0.778	38 8%	4 1%	0.823
T2b (yes)	123	28.5	24 28%	1.14	228 32	72 10%	1.14
GENDER	80	18.5	20 4%	0.778	22 32	20 6%	0.778
Stage	5	51	10%	11 2%	46 9%	16 3%	5
T2b	115	24.5	22 4%	0.88	188 49	49 10%	0.88

\*Low and high-exp were determined by 4-means clusterization. \*\*Based on Pearson's  $\chi^2$  test.



# More visualizations



# Filter pubmed by impact factors

SJR (1,000)

sjr:>10

sjr:markdown

Scimago Journal & Country Rank (SJR)

http://www.scimagojr.com/

Pubmed (168173)

sjr:multi

United States  
United Kingdom  
Netherlands  
Germany  
Sweden

sjr:table\_by\_title

title.keyword: Descending

Count

title.keyword: Descending	Count
Annual Review of Biochemistry	18
Annual Review of Cell and Developmental Biology	18
Annual Review of Genetics	18
Annual Review of Immunology	18
Annual Review of Neuroscience	18
Cell	18
Chemical Reviews	18
Genes and Development	18
Immunity	18
Molecular Cell	18

Export: Raw Formatted

1,000 Count  
16.193 Average SJR  
10,299.874 Average cites 3 years  
9,582.665 Total references

sjr:enhancedTable

title	citable_doc	cites_report_doc_2_years	country	issn	sjr	sjr_best_quartile	total_cites_3_years	total_docs_3_years	total_docs_years	total_ref	type
Astrophysical Journal Letters	7	72.75	United Kingdom	20418213, 20418205	62,429	Q1	493	7	5	350	journal
Annual Review of Cell and Developmental Biology	60	26.66									
Current Opinion in Cell Biology	303	23.29									
Neuron	752	15.47									
Annual Review of Astronomy and Astrophysics	47	13.23									

Filter...  
... related to (101) from SJR

pubmed:multi\_by\_journal

Nature  
Journal of the American Chemical Society  
Cell  
Neuron  
Immunity  
National vital statistics Reports  
Rhinology Supplement  
NIH consensus and statements  
Science

pubmed:metrics

168,173 Number of references  
9 Number of journals

pubmed:metabolites  
HMDB\_metabolite (235)  
HMDB\_protein (3131)  
Pubtator (103762) SJR (101)

pubmed:tag\_by\_meshLabel

Transcription, Genetic  
Signal Transduction  
Amino Acid Sequence  
United States  
Aged  
Humans  
Middle Aged  
Mice  
DNA  
Base Sequence Research  
Rats  
Male  
Female  
Adult Cells, Culture  
Molecular Sequence Data  
Mutation  
Transcription Factors  
Time Factors  
RNA, Messenger  
Gene Expression Regulation

pubmed:table\_by\_AuthorLastName

Author name	Count
Wang	2,610
Chen	2,123
Lee	2,046
Li	1,995
Zhang	1,924
Smith	1,763
Liu	1,637
Kim	1,416
Wu	1,201

pubmed:tag\_by\_chemicalName

Macromolecular Substances  
Recombinant Fusion Proteins  
Saccharomyces cerevisiae Proteins  
Membrane Proteins  
Carrier Proteins  
RNA, Messenger  
Peptides  
Proteins  
Transcription Factors  
Nuclear Proteins  
Adenosine Triphosphate  
Bacterial Proteins  
Drosophila Proteins  
RNA  
DNA  
DNA-Binding Proteins  
Proteins  
RNA  
DNA  
Messenger  
Bacterial Proteins  
Drosophila Proteins  
Nerve Tissue Proteins  
Receptors, Cell Surface

pubmed:enhancedTable

PMID	Title	Year	Journal
-	Novel DNA polymerases offer clues to the molecular basis of mutagenesis.	1,999	Cell
-	Epigenetic spreading of the Drosophila dosage compensation complex from roX genes into flanking chromatin.	1,999	Cell
-	Membrane polarizing controlling the sorting and diffusion of membrane components.	1,000	Neuron

pubmed:table\_by\_collectiveNames

Collective name  
American Heart Association  
Department of Biochemistry and Molecular Biology, Harvard University, Cambridge, Massachusetts 02138, 45

pubmed:table\_by\_affiliation

Affiliation  
Department of Biochemistry and Molecular Biology, Harvard University, Cambridge, Massachusetts 02138, 87

98

# A little overview of Kibio

**Biomedical literature**

- pubmed 27.8m
- SJR 1k
- pubtator 21.3m
- HMDB
- t t3DB
- G GTex
- M MeSH
- descriptor 374.1k
- pharmacological 6.2k
- qualifier 1.2k
- supplemental 859.3k
- SMPDB
- Others
- disgenet\_disease\_gene 1.5m
- bcac 35.4m
- snpdb 325m
- drugbank 8.3k
- gencode 1.9m
- Omim 25.1k
- Samples
- er\_transcript\_abundance 348
- triceps 835.3k

**pubmed (27,766,141)**

Filter...

**pubmed:bar\_by\_date**

References per year

**pubmed:multi\_by\_journal**

27,766,141 Number of references

30,662 Number of journals

**pubmed:relations**

- HMDB\_metabolite (2121) (\*)
- HMDB\_protein (2305) (\*)
- Pubtator (3842503) (\*)
- SJR (101)

**pubmed:tag\_by\_meshLabel**

Follow-Up Studies  
Aged, 80 and over United States

Infant, Newborn Infant Animals Male Cells, Cultured  
Child, Preschool Rats Female Adult Child Pregnancy  
Middle Aged Humans Aged Young Adult  
Risk Factors Adolescent Mice Time Factors  
Treatment Outcome Retrospective Studies  
Molecular Sequence Data

**pubmed:table\_by\_AuthorLastName**

Author name	Count
Wang	642,966
Li	525,719
Zhang	511,305
Chen	469,437
Liu	437,649
Lee	387,326
Kim	307,483
Yang	292,313
Wu	249,136

**pubmed:tag\_by\_chemicalName**

Recombinant Proteins Biomarkers  
Blood Glucose Anti-Bacterial Agents  
Oxygen RNA, Messenger  
Transcription Factors Proteins DNA Calcium  
Peptides Water Antineoplastic Agents  
Antibodies, Monoclonal Membrane Proteins  
Amino Acids Glucose Bacterial Proteins Insulin DNA-Binding Proteins

**pubmed:enhancedTable**

PMID	Title	Year	Journal
3874232	Physicochemical and functional properties of murine B cell-derived B cell growth factor II (WEHI-231-BCGF-II).	1,985	Journal of immunology (Baltimore, Md.: 1950)
-	An analysis of monoclonal T cell and antibody recognition sites on Ia molecules.	1,985	The Journal of investigative dermatology
-	Prolonged survival in vivo of unprimed B cells responsive to a T-independent antigen.	1,985	The Journal of experimental medicine

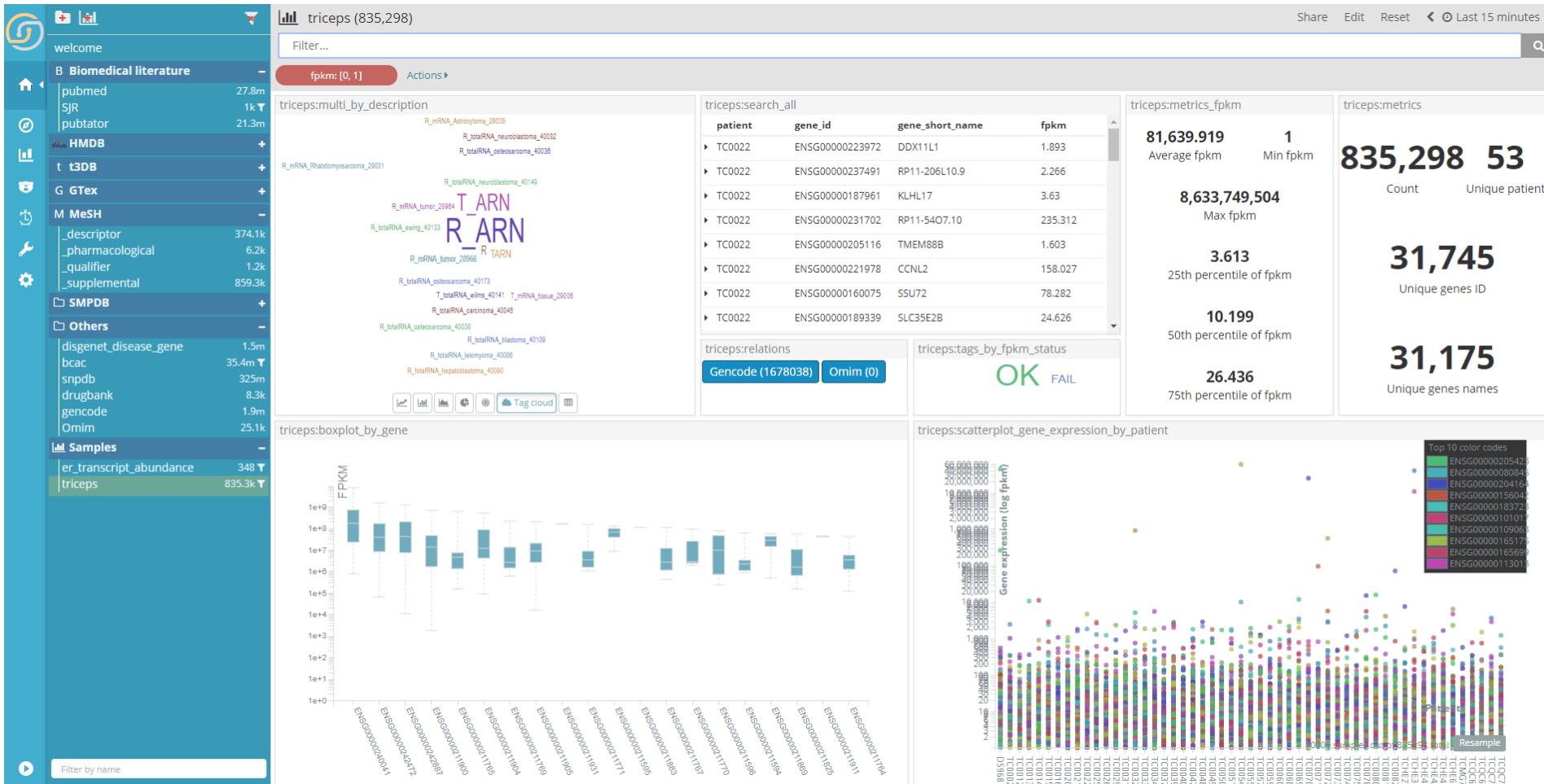
**pubmed:table\_by\_collectiveNames**

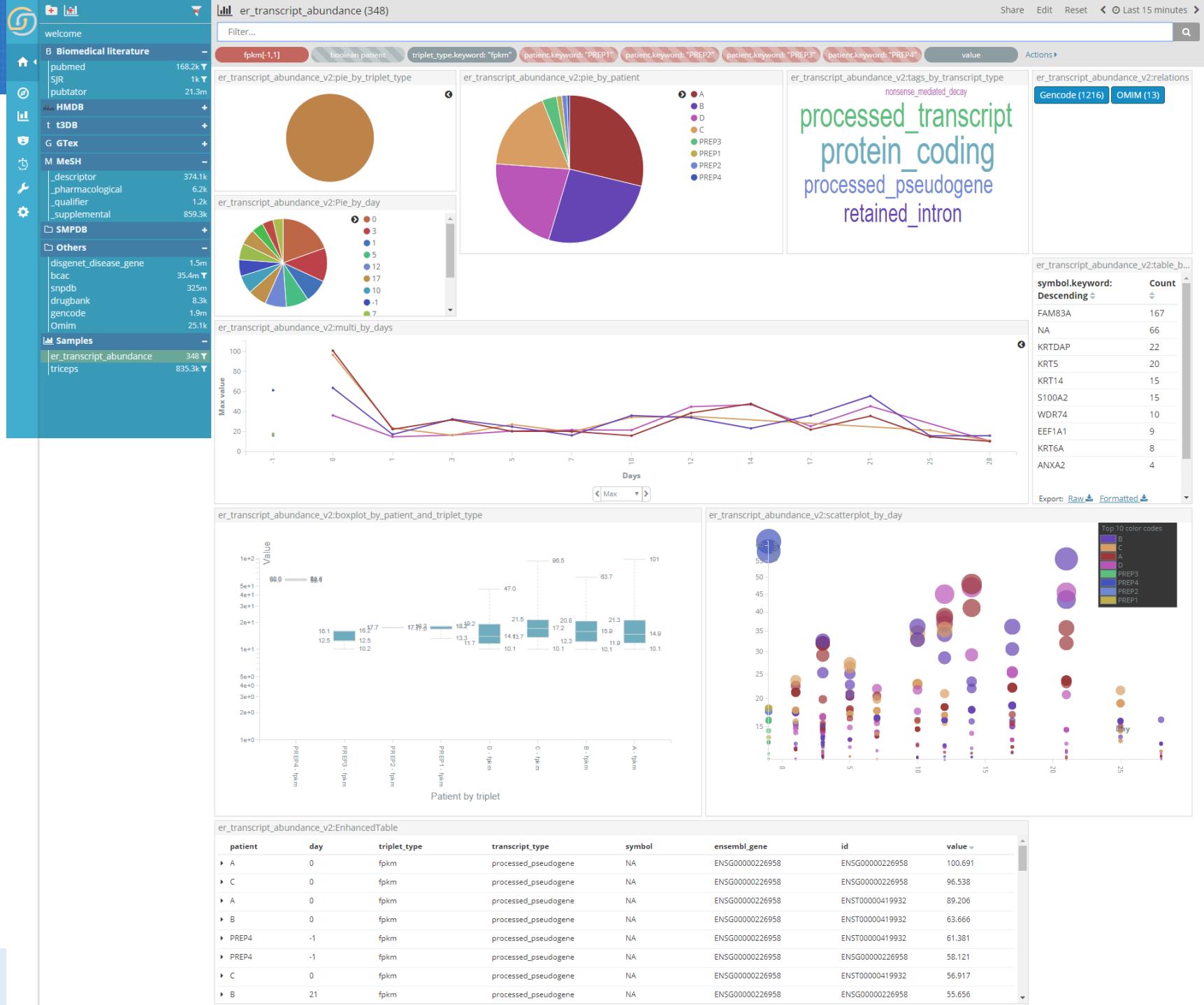
Collective name	Count
Centers for Disease Control and Prevention (CDC)	4,691
Centers for Disease Control (CDC)	2,253

**pubmed:table\_by\_affiliation**

Affiliation	Count
The BMJ	1,330
Laboratories of The Rockefeller Institute for Medical Research.	1,272
Department of Psychology.	1,218

# Transcriptomics data





# Conclusions

- Fast searches within linked databases
- Better exploration
- Intuitive portal
- Import new data
- Future implementations:
  - Analytical tools (PCA, cohorts comparison)
  - Biomarker discovery tools

# Global view

