

Omics data analysis for high-throughput phenotyping

Etienne Thévenot *et al.*

CEA, LIST (Saclay, France)

Laboratory for Data Analysis and Systems' Intelligence
MetaboHUB

etienne.thevenot@cea.fr

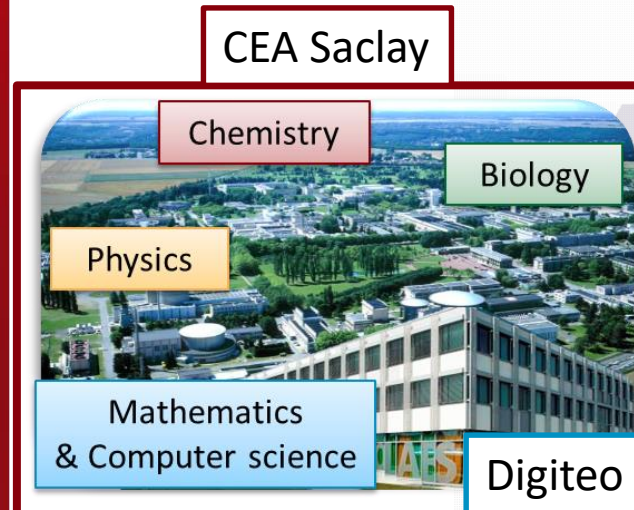
<http://etiennethevenot.pagesperso-orange.fr/>

FROM RESEARCH TO INDUSTRY

cea

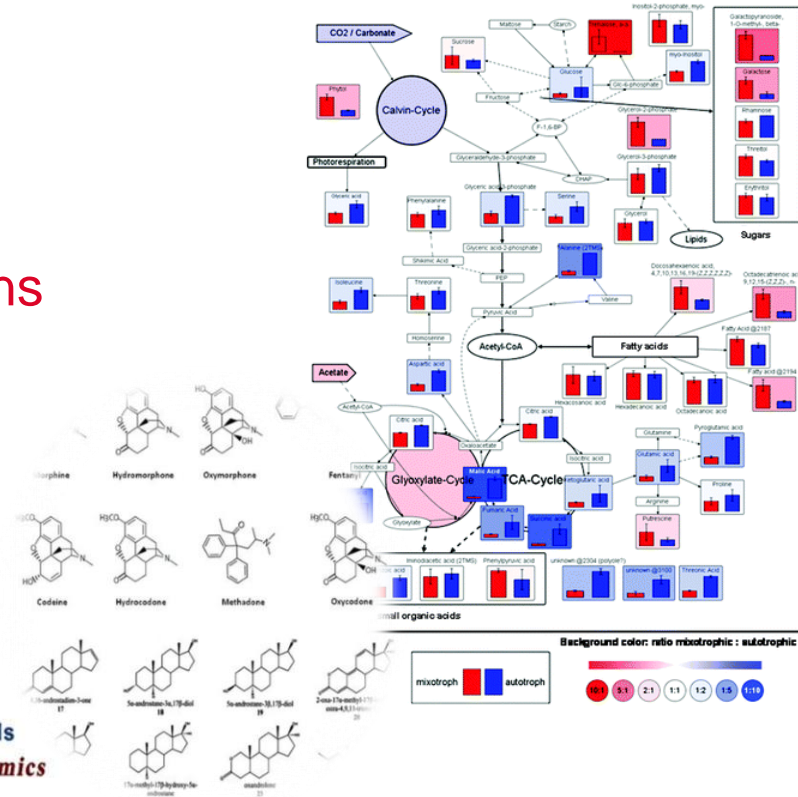
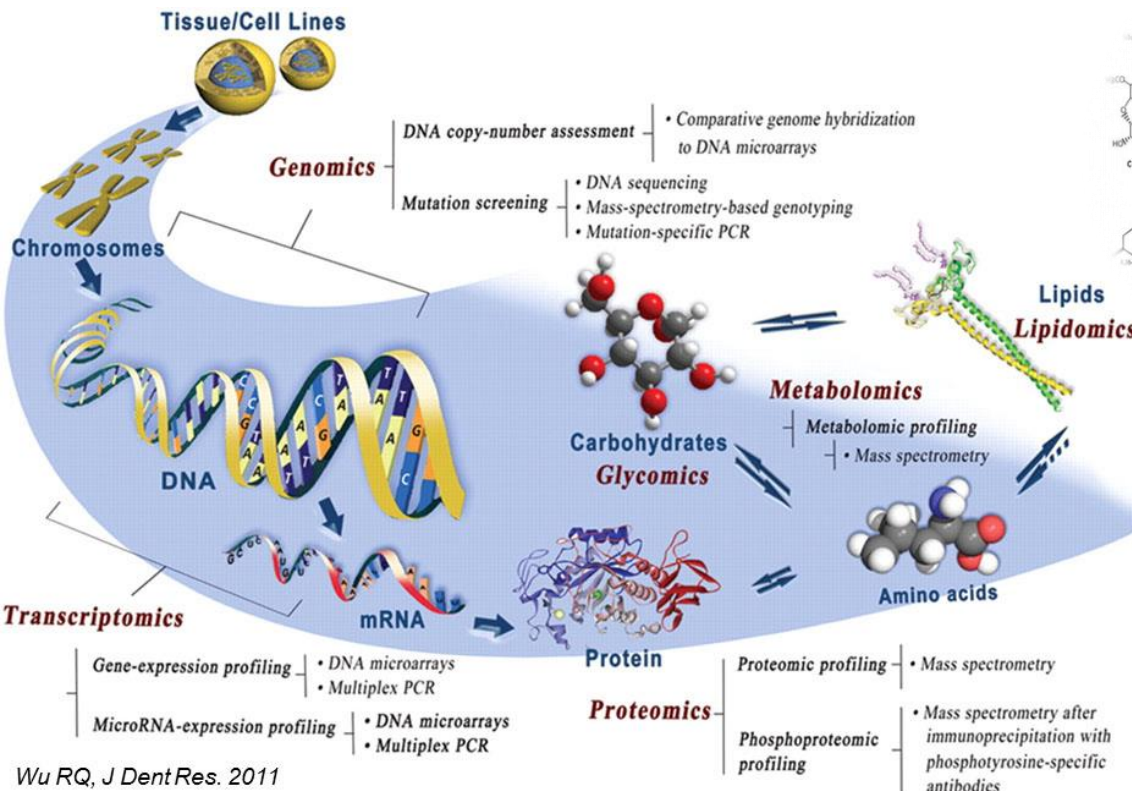


www.cea.fr



➤ **Metabolomics**

- omics science
- dedicated to small molecules (< 1kDa)
- involved in metabolic chemical reactions



FUN MOOC La métabolomique : enjeux technologiques et scientifiques

Metabolomics successes for clinical biomarkers

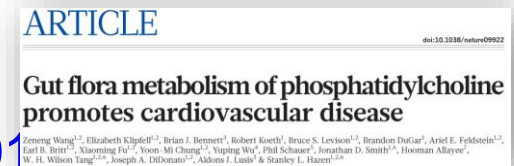
■ **Cancer** and oncometabolites: [Yang et al., 2013](#)



■ **Memory** impairment and phospholipids: [Mapstone et al., 2014](#)



■ **Cardiovascular** disease and TMAO: [Wang et al., 2011](#)



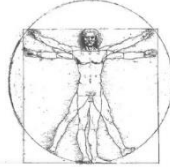
■ **Diabetes** and amino acids: [Wang et al., 2011](#)



■ **Review**: [Wishart et al., 2016](#)



cea Metabolomics workflow



Biological question → Design of Experiment



Data integration

Analytical Chemistry



NMR



LC-MS
GC-MS

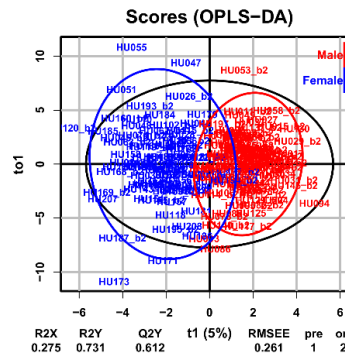
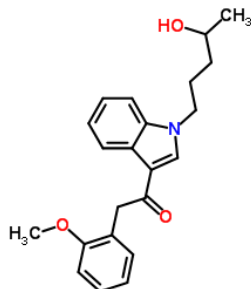
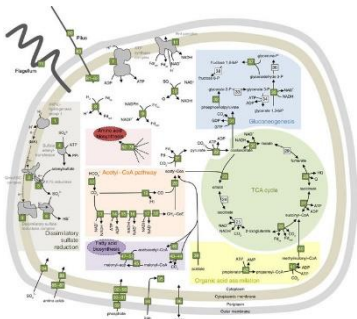
Network analysis

Preprocessing

Chemical identification

Statistics

mz	rt	Db_015	...	Db_068
75.0322	41.28	22162	...	48575
75.0441	174.83	1371	...	820
75.0634	56.23	49111	...	91769
...
999.6653	844.61	571	...	636
999.6759	844.61	711	...	665
999.6865	844.61	698	...	612



Biological question Design of Experiment

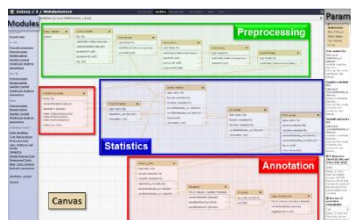
Data integration



Analytical Chemistry

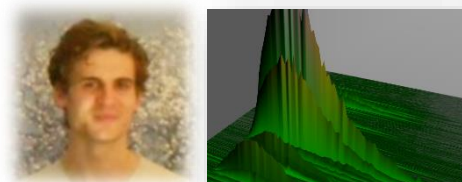


Workflow4Metabolomics



Network analysis

Preprocessing



proFIA

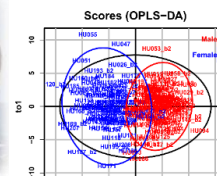
Gene Expression
proFIA: A data preprocessing workflow for Flow Injection Analysis coupled to High-Resolution Mass Spectrometry
Alexis Delabrière^{1,*}, Ulli M. Hohenester², Benoit Colsch², Christophe Junot², François Fenaille² and Etienne A. Thévenot^{1,*}

Chemical identification

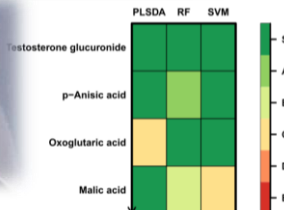
Statistics



ropls

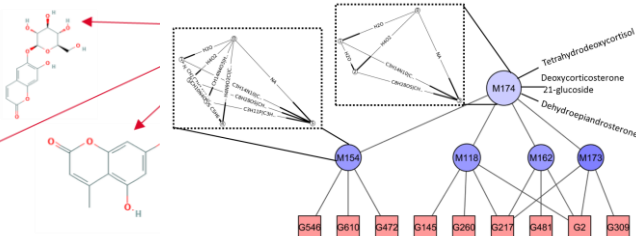


biosigner



MineMS2

rbiodb

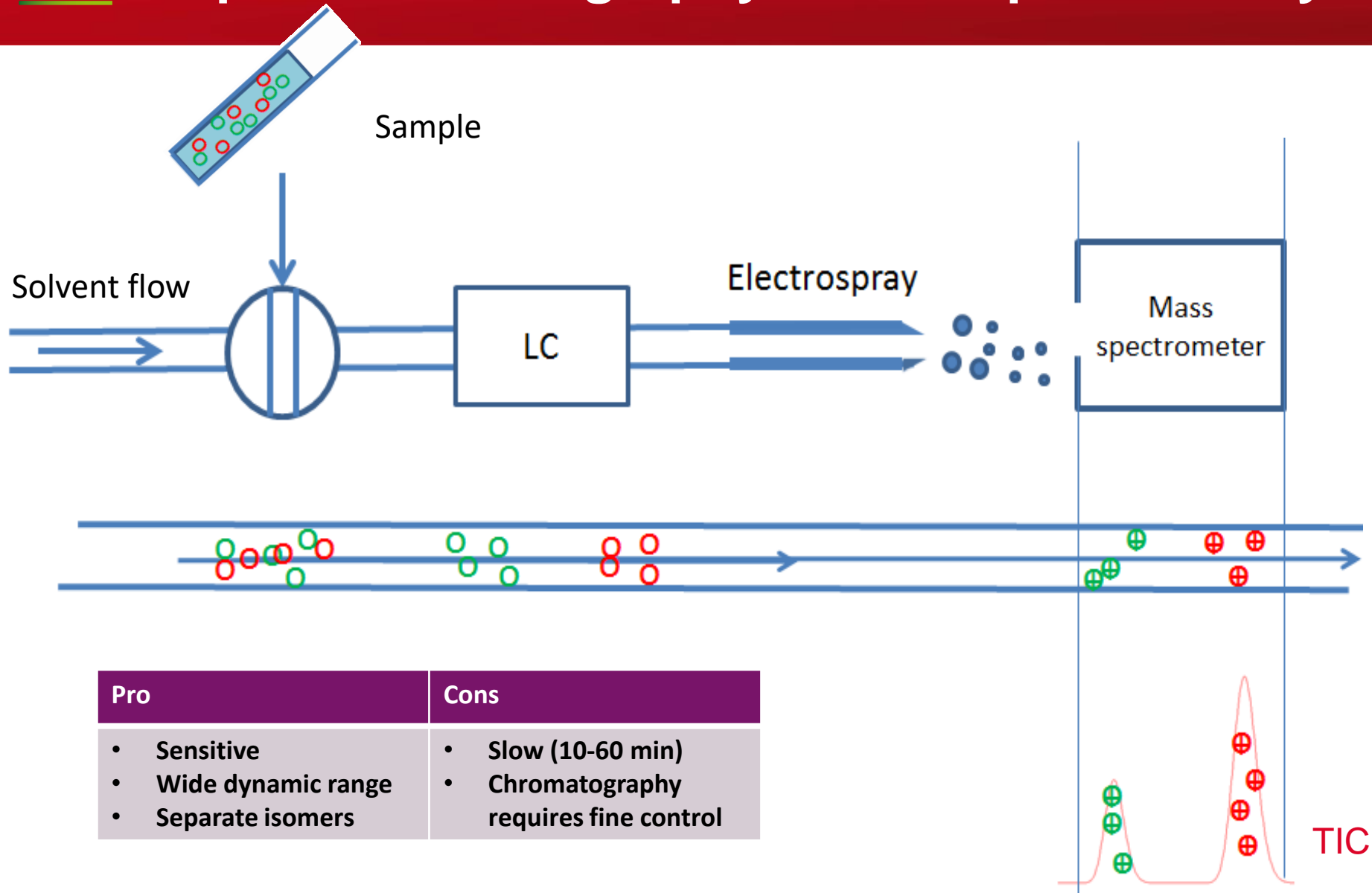


Journal of proteome research
Analysis of the Human Adult Urinary Metabolome Variations with Age, Body Mass Index, and Gender by Implementing a Comprehensive Workflow for Univariate and OPLS Statistical Analyses
Etienne A. Thévenot,^{1,*} Aurélie Roux,^{2,†} Ying Xu,² Eric Ezan,² and Christophe Junot^{2,*}

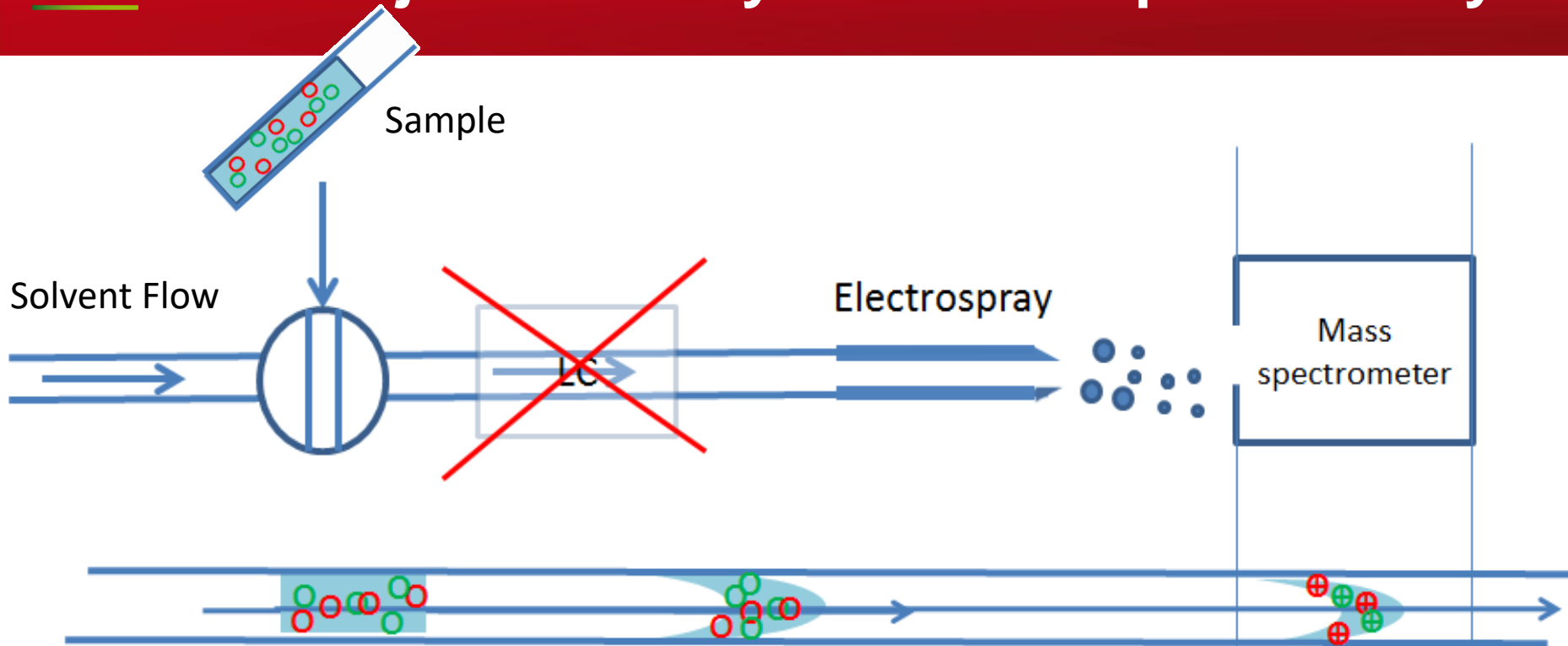
frontiers in Molecular Biosciences
biosigner: A New Method for the Discovery of Significant Molecular Signatures from Omics Data
Philippa Rinaudo¹, Samia Boudah¹, Christophe Junot² and Etienne A. Thévenot^{1,*}

- **Data preprocessing**
 - **Flow injection analysis**

Liquid chromatography - mass spectrometry

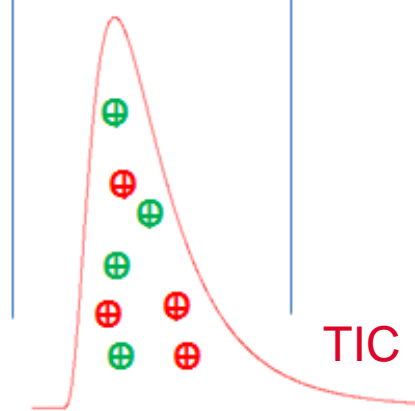


Flow injection analysis - mass spectrometry



	Pro	Cons
LC-MS	<ul style="list-style-type: none"> • Sensitive • Wide dynamic range 	<ul style="list-style-type: none"> • Slow (10-60 min) • Expensive and difficult chromatography
FIA-MS	<ul style="list-style-type: none"> • Fast (1-3 min) 	<ul style="list-style-type: none"> • Matrix effect • Isomers not separated

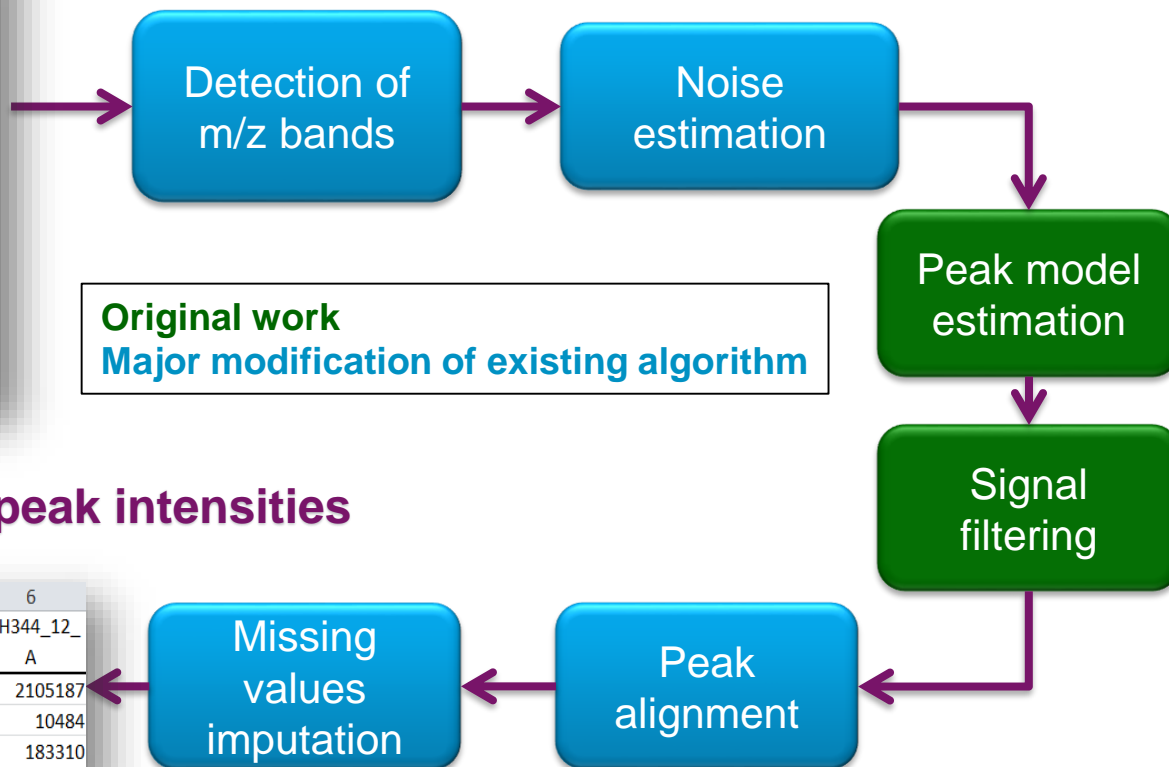
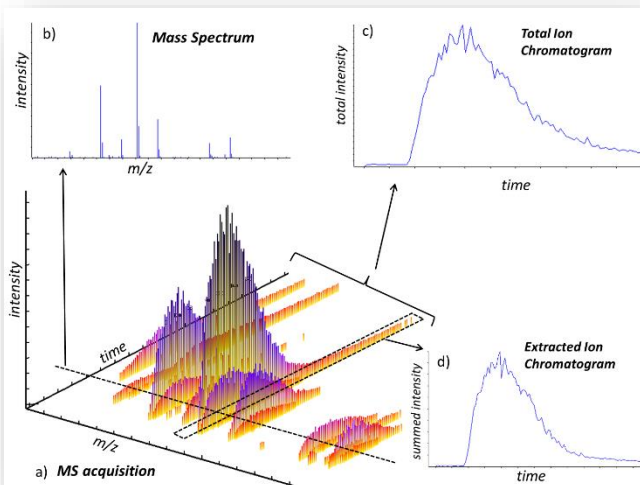
FIA is adapted to high-throughput screening





Raw files

Alexis Delabrière

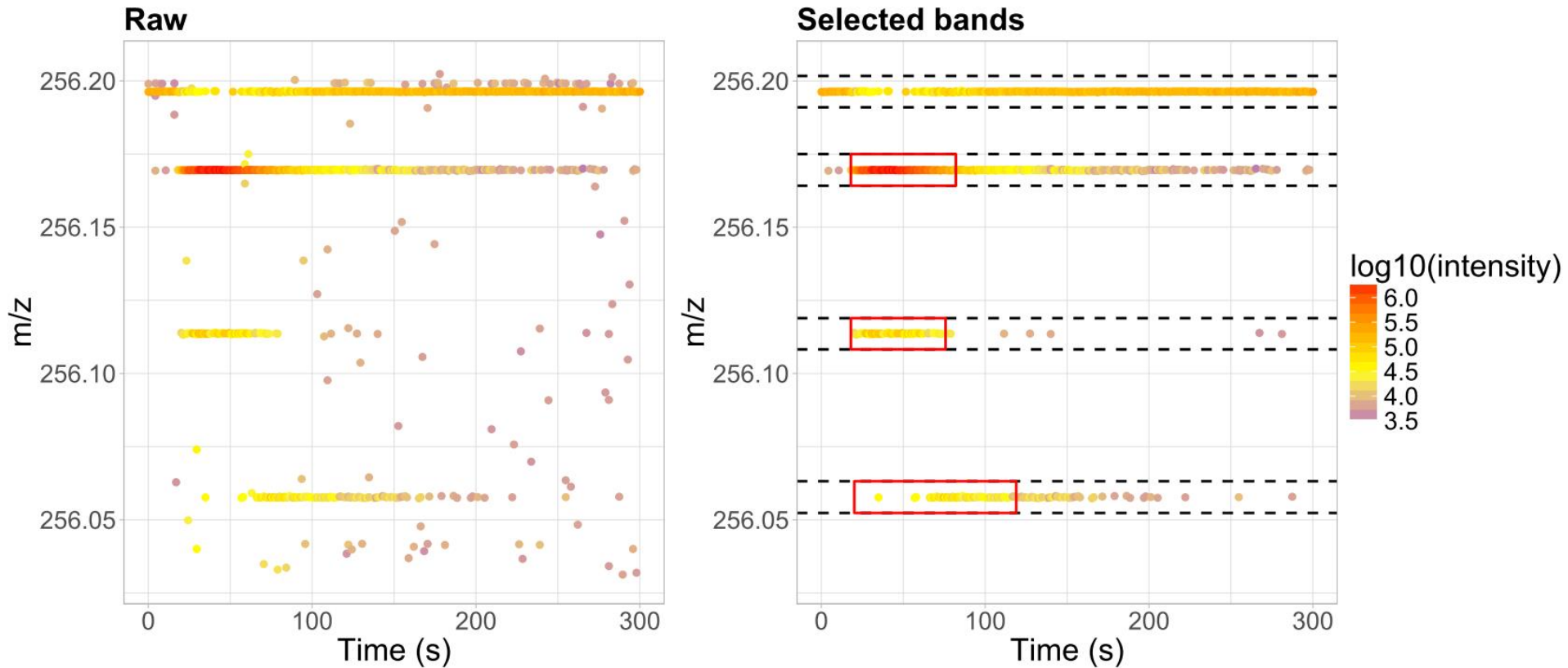


Variable by sample table of peak intensities

	1	2	3	4	5	6
1	mzMed	mzMin	mzMax	meanSolvent	corMean	UIH344_12_A
2	163.02777	163.02767	163.02782	845186.744	0.42770334	2105187
3	163.03893	163.03889	163.039	0	0.71292516	10484
4	163.11577	163.11562	163.11586	0	0.54386442	183310
5	164.02923	164.02919	164.0294	0	0.6646757	119575

Delabriere *et al.* (2017). *proFIA*: A data preprocessing workflow for Flow Injection Analysis coupled to High-Resolution Mass Spectrometry. *Bioinformatics*. 33:3767-3775.

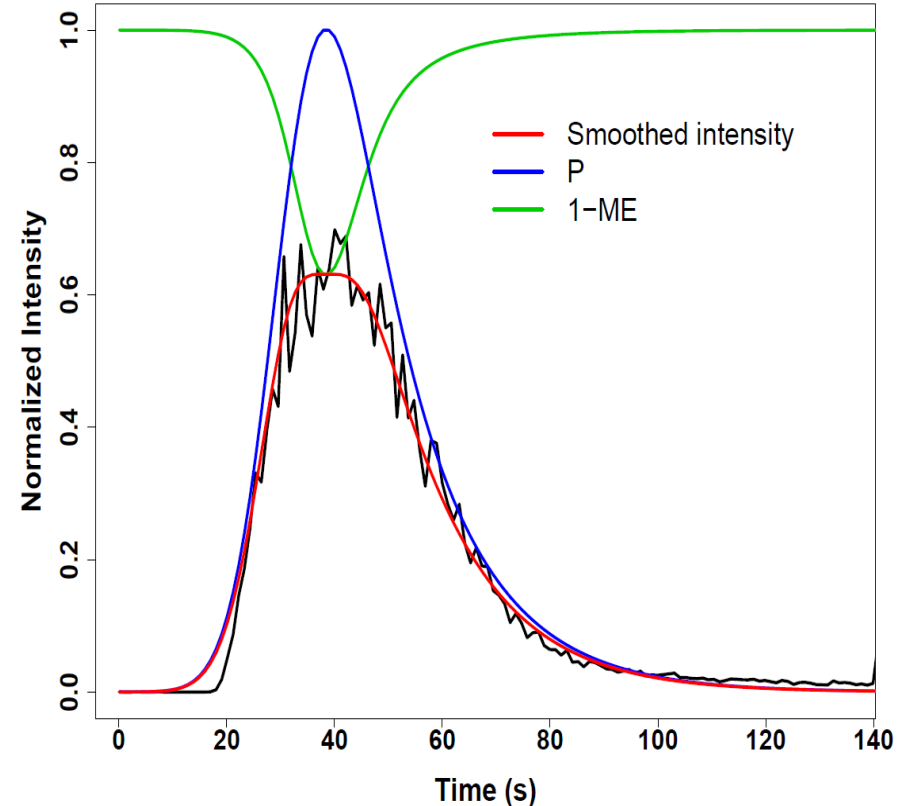
Detection of m/z bands



➤ We proposed a model based on Kolev (1994) and Nanita (2012)

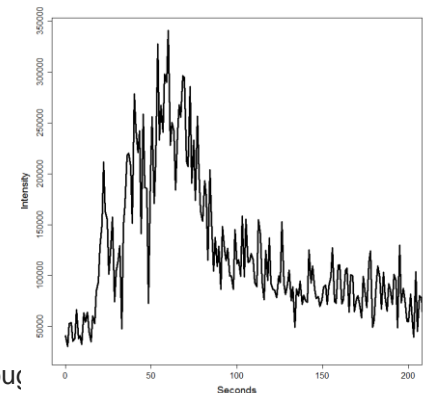
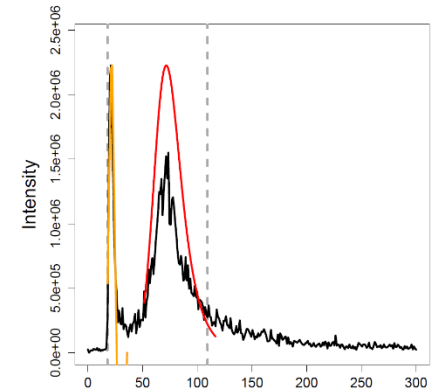
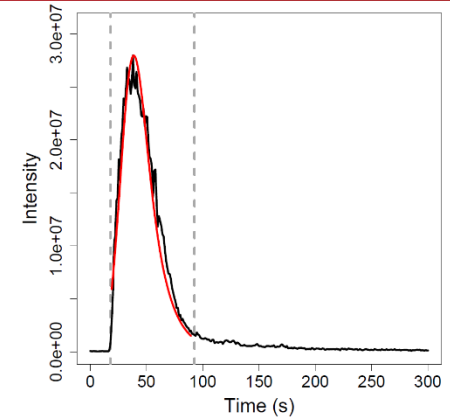
With :

- I_A the observed intensity
- k_A a constant specific to the molecule
- P is exponentially modified gaussian
- ME_A is a second order exponential
- B_A is the baseline constant for analyte
- ϵ is the heteroscedastic noise



➤ Intense peaks without baseline are selected and a regression is performed leading to a peak model P

- This peak model is used to perform matching filtration on the signal
- The match can be extended if a second maximal is found on the filter. If not, a triangular filter is used for coarser grain
- A statistical test has been developed to discard signals too close to the baseline



cea Application to metabolomics data

➤ Dataset:

- plasma sample spiked with 40 molecules at 6 concentrations

➤ Running time:

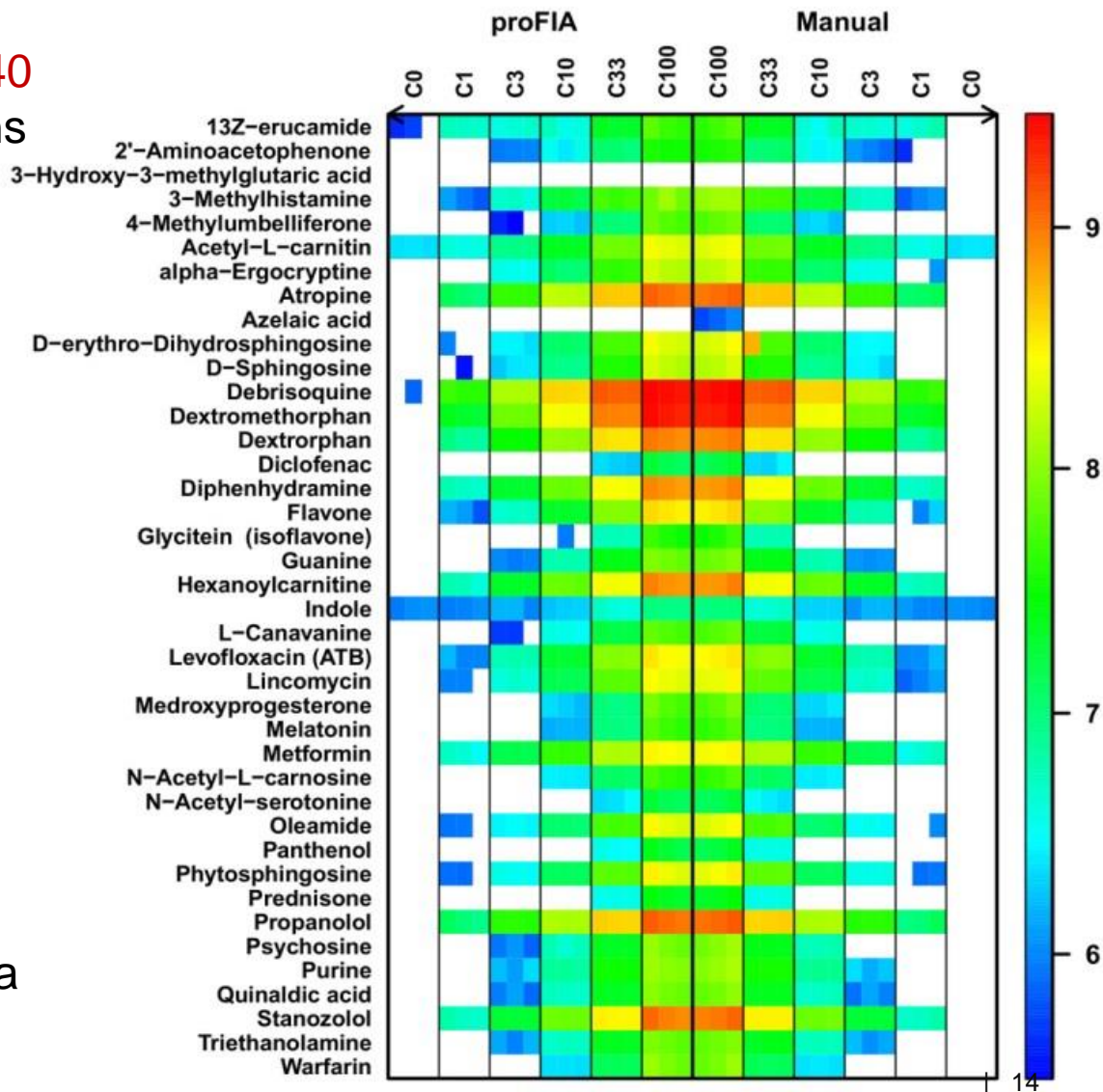
- < 15 s per file

➤ Comparison with manual integration:

- precision of 0.96
- recall of 0.98
- mean intensity error < 5%

➤ Annotation:

- 211 signals out of 1082 had a unique match on HMDB



cea The *proFIA* software

- R package: Bioconductor
([DOI:10.18129/B9.bioc.proFIA](https://doi.org/10.18129/B9.bioc.proFIA))



- Galaxy tool: Toolshed, Workflow4Metabolomics, PhenoMeNal



- Publication: *Bioinformatics*
([DOI:10.1093/bioinformatics/btx458](https://doi.org/10.1093/bioinformatics/btx458))

Bioinformatics

doi:10.1093/bioinformatics/xxxxx

Advance Access Publication Date: Day Month Year

Manuscript Category

Gene Expression

***proFIA*: A data preprocessing workflow for Flow Injection Analysis coupled to High-Resolution Mass Spectrometry**

Alexis Delabrière^{1,*}, Ulli M. Hohenester², Benoit Colsch², Christophe Junot², François Fenaille² and Etienne A. Thévenot^{1,*}

LC-MS

Preprocessing

xcms.xcmsSet Filtration and Peak Identification using xcmsSet function from xcms R package to preprocess LC/MS data for relative quantification and statistical analysis

xcms.xcmsSet Merger Merge xcms.xcmsSet xset in one to be used by group

xcms.group Group peaks together across samples using overlapping m/z bins and calculation of smoothed peak distributions in chromatographic time.

xcms.retcor Retention Time Correction using retcor function from xcms R package

xcms.fillPeaks Integrate a sample's signal in regions where peak groups are not represented to create new peaks in missing areas

xcms.summary Create a summary of XCMS analysis

CAMERA.annotate CAMERA annotate function. Returns annotation results (isotope peaks, adducts and fragments) and a diffreport if more than one condition.

CAMERA.combinexsAnnos Wrapper function for the combinexsAnnos CAMERA function. Returns a dataframe with recalculated annotations.

proFIA Preprocessing of FIA-HRMS data
metaMS.wanGC GC-MS data preprocessing using metaMS package

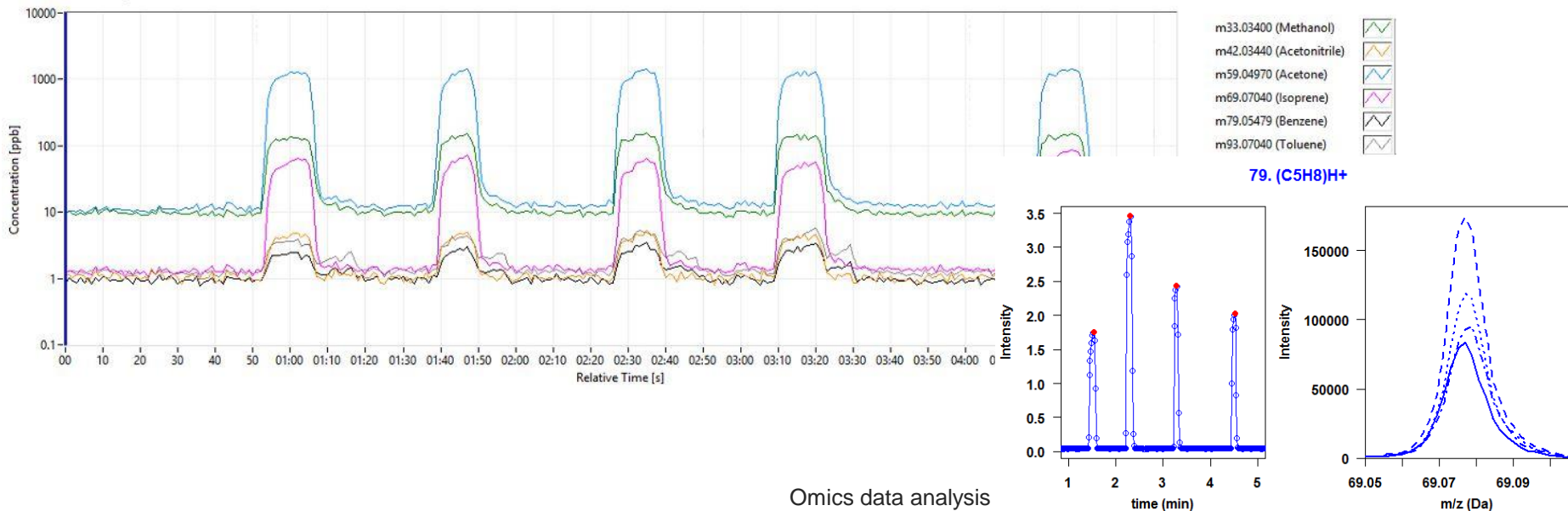
Breathomics: real time analysis of exhaled air in disease and response to treatments

- Objective: comprehensive analysis of metabolism-derived Volatile Organic Compounds (VOCs)
- Technology: PTR-TOF-MS at the patient bedside (Foch Hospital)
- Project: develop innovative algorithms and software environment for the processing of real-time analysis of VOCs in exhaled air



HOPITAL
FOCH

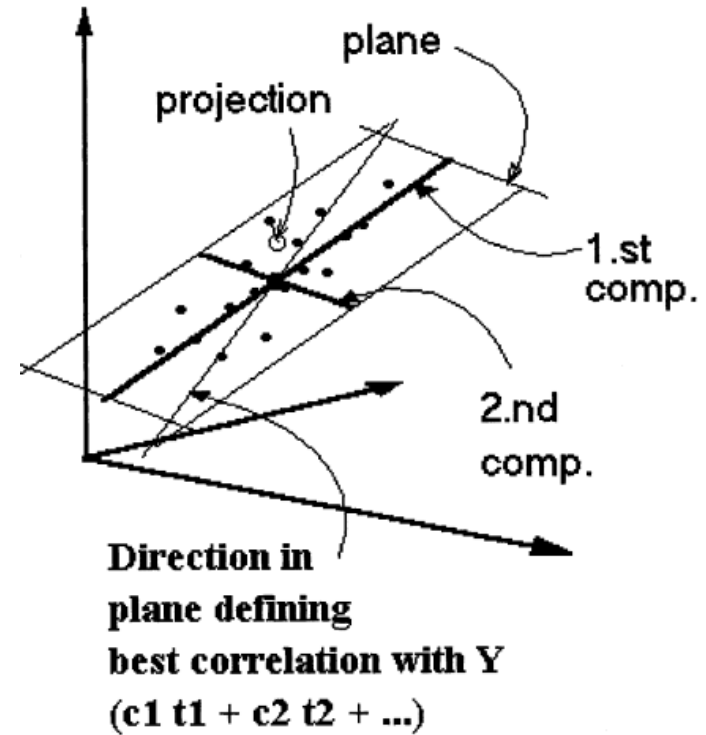
S. Grassin-Delyle



➤ **Statistical analysis**

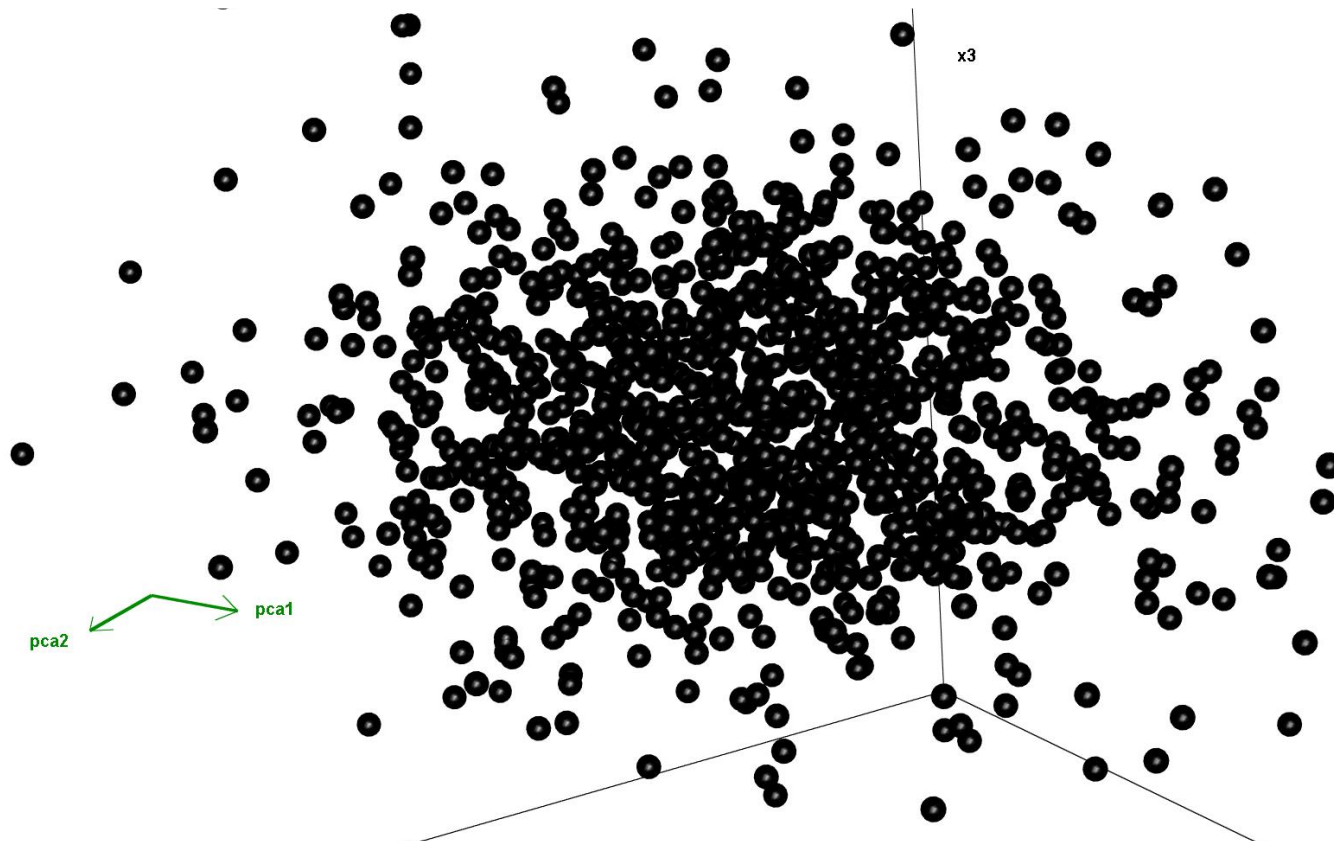
- **Orthogonal Partial Least Squares**

- Multivariate regression approach
- Handle data sets
 - of high dimension ($n < p$)
 - correlated variables
 - including missing values
- Based on latent variables
 - maximizing covariance with the response Y
- Developed by Wold H. and S.
- Can be used for classification (PLS-DA)



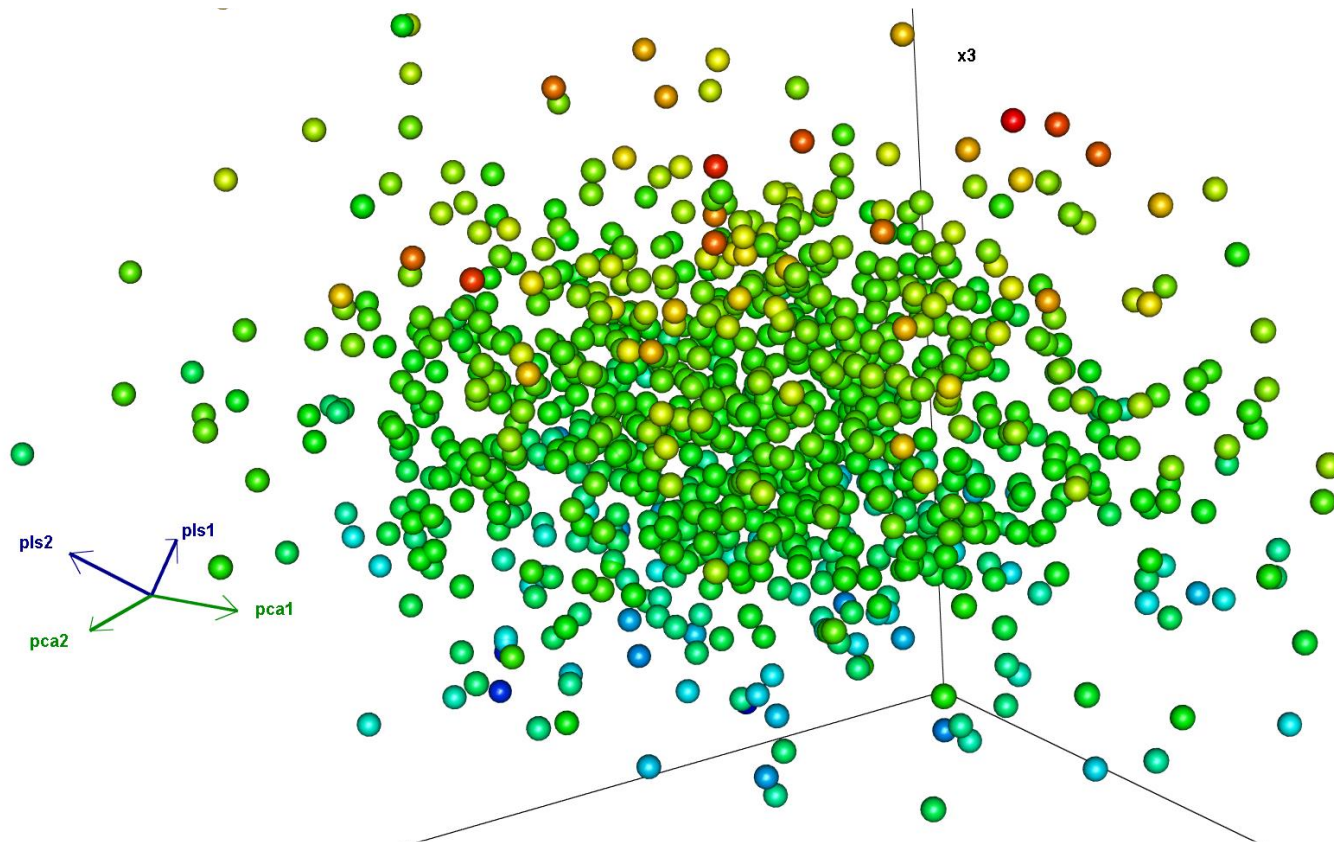
[Wold et al. \(2001\). PLS-regression: a basic tool of chemometrics. Chemometrics and Intelligent Laboratory Systems 58, 109–130.](#)

➤ PCA finds the directions of maximum variance

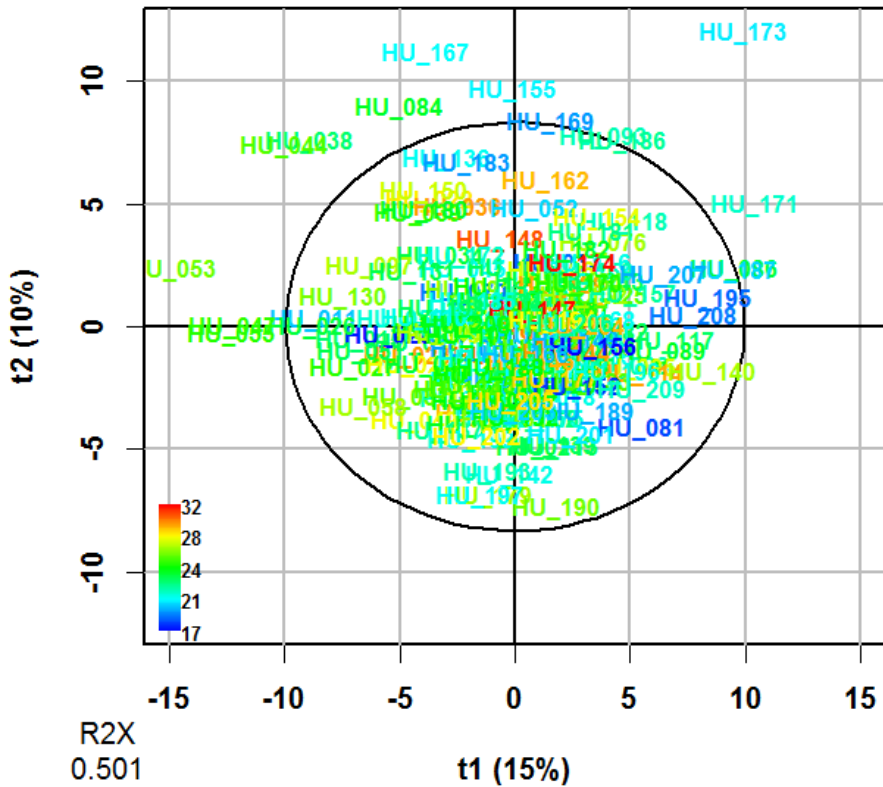


cea PLS vs PCA: score plots

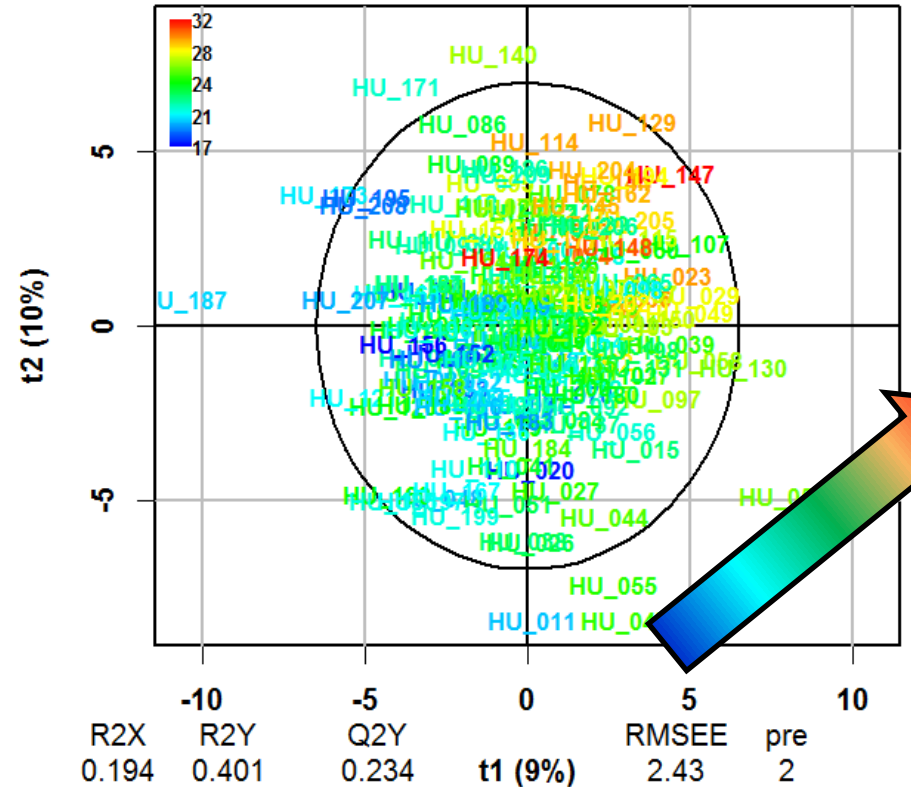
➤ PLS includes the labels into the model



Scores (PCA)



Scores (PLS)



ropls package: R implementation of the (O)PLS(-DA) modeling algorithms

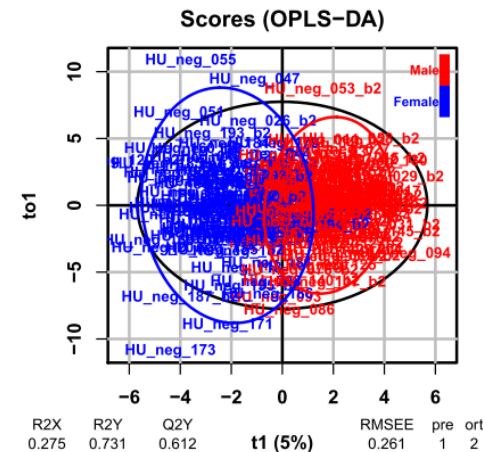
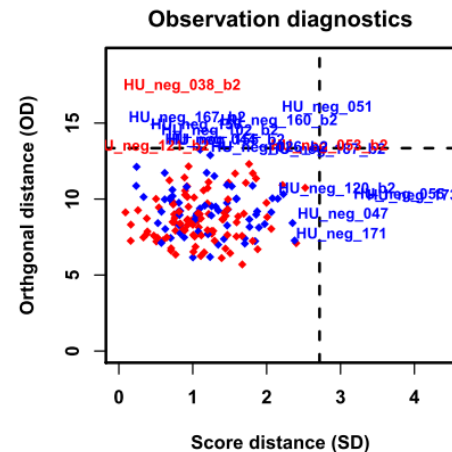
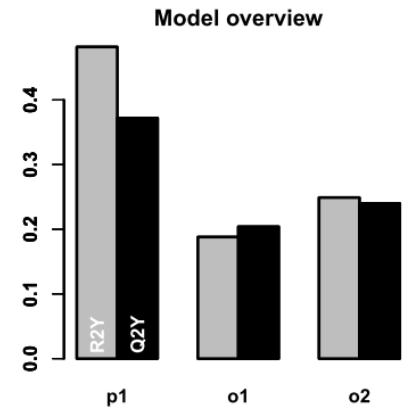
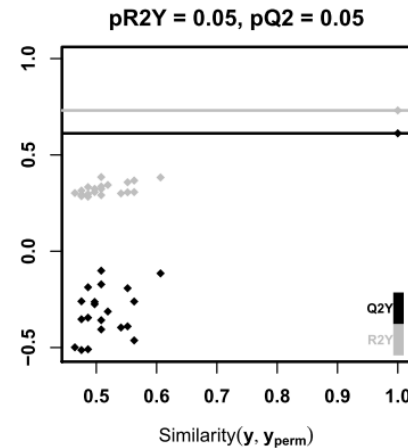


➤ Full diagnostics

- outliers
- permutation testing

➤ Full numerical and graphical results

- R2X, R2Y, Q2Y
- VIPs

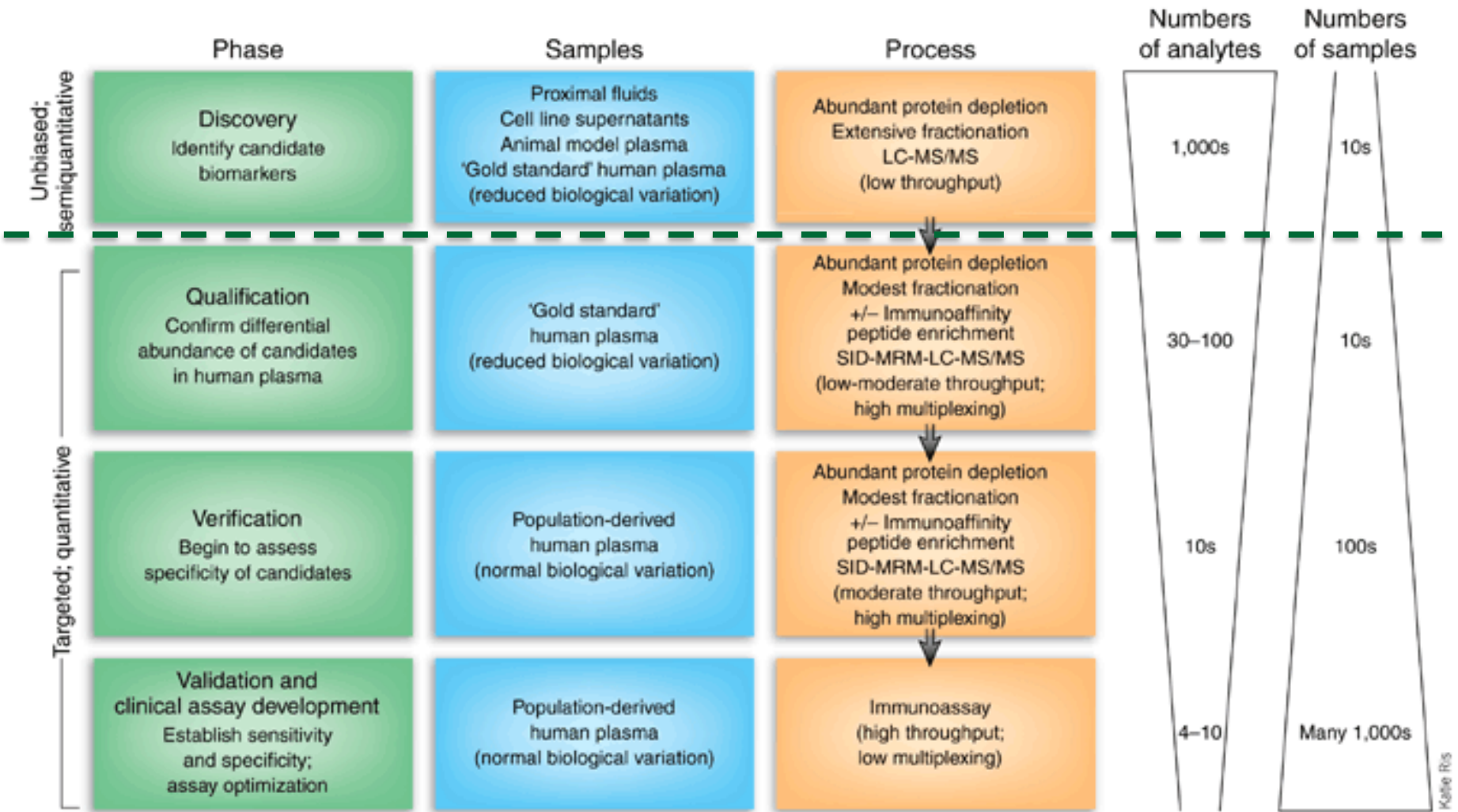


[Thévenot *et al.* \(2017\). Analysis of the human adult urinary metabolome variations with age, body mass index and gender by implementing a comprehensive workflow for univariate and OPLS statistical analyses. *J. Proteome Res.* 14:3322-3335.](#)

➤ **Statistical analysis**

- **Feature selection**

Feature selection: from biomarker discovery to clinical diagnostics



[Rifai et al. \(2006\). Protein biomarker discovery and validation: the long and uncertain path to clinical utility. Nat. Biotechnol. 24:971-983.](#)

- Restrict the list of candidates before the subsequent validation phases
- Facilitates interpretation
- Limit the risk of overfitting
- Stabilize the prediction

Feature selection: challenges

- Testing all combination of features is not computationally tractable
 - efficient search path

- Prediction performance
 - sensitivity, selectivity

- Stability
 - reproducibility

- Relevance
 - selection criterion



Feature selection: approaches

➤ filter (threshold criterion)

- e.g., t-test



fast

threshold?

➤ wrapper (iterative selection)

- e.g., SVM RFE, Genetic Algorithm

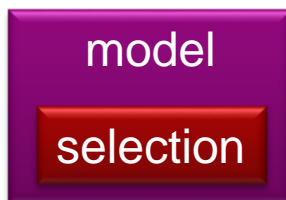


interaction
with classifier

computation
intensive

➤ embedded (penalization constraint)

- e.g., Lasso, Elastic Net



fast

stability

➤ **Statistical analysis**

- The *biosigner* approach



Philippe Rinaudo

- Objective: select only features which significantly contribute to the performance of the classifier
- Method: features are significant if the prediction accuracy decreases after random permutation of their values for in test samples
- Algorithm:
 1. generate k train/test subset by resampling
 2. build the models and rank the variables
 3. find the largest non-significant feature subset (half-interval search)
 4. repeat steps (1-3) on the dataset restricted to the significant features until the selection is stable (all features are significant)

Rinaudo et al. (2016). *biosigner*: a new method for the discovery of significant molecular signatures from omics data. *Front. Mol. Biosci.* 3.

1.1 Generate k subsets (bootstrap resampling)

response (y)

dataset(X)



1.1 Generate k subsets (bootstrap resampling)

$test_k$



$train_k$



1.2 Train F_k models (e.g. PLS-DA)

$$y_k = F_k (X_k)$$



1.3 Find the largest subset of non-significant features

b) Find the feature f of lowest rank such that the subset of all features of higher ranks is not significant:

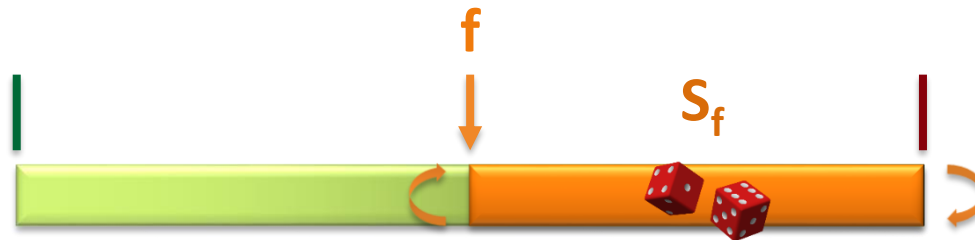
i) set f to the feature of mean rank



1.3 Find the largest subset of non-significant features

b) Find the feature f of lowest rank such that the subset of all features of higher ranks is not significant:

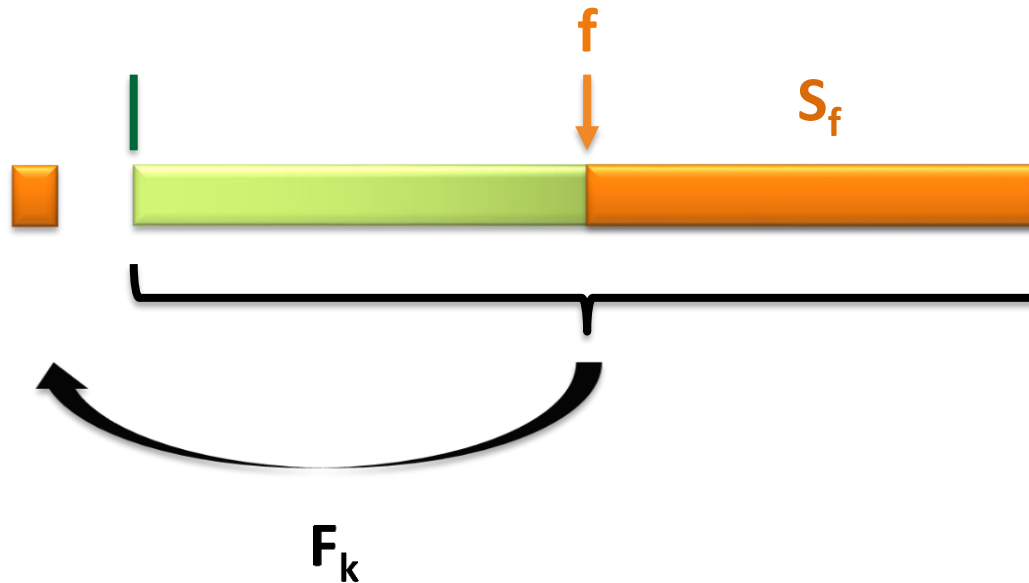
ii) permute in the test set all features of higher rank (ie features in S_f)



1.3 Find the largest subset of non-significant features

b) Find the feature f of lowest rank such that the subset of all features of higher ranks is not significant:

iii) compare the accuracies of the predictions after permutation



1.3 Find the largest subset of non-significant features

b) Find the feature f of lowest rank such that the subset of all features of higher ranks is not significant:

iii-a) accuracy \rightarrow or \nearrow

$\Rightarrow S_f$ does not contain significant features



1.3 Find the largest subset of non-significant features

b) Find the feature f of lowest rank such that the subset of all features of higher ranks is not significant:

iii-a) accuracy \rightarrow or \nearrow

=> shift f upward to the mean rank of significant features

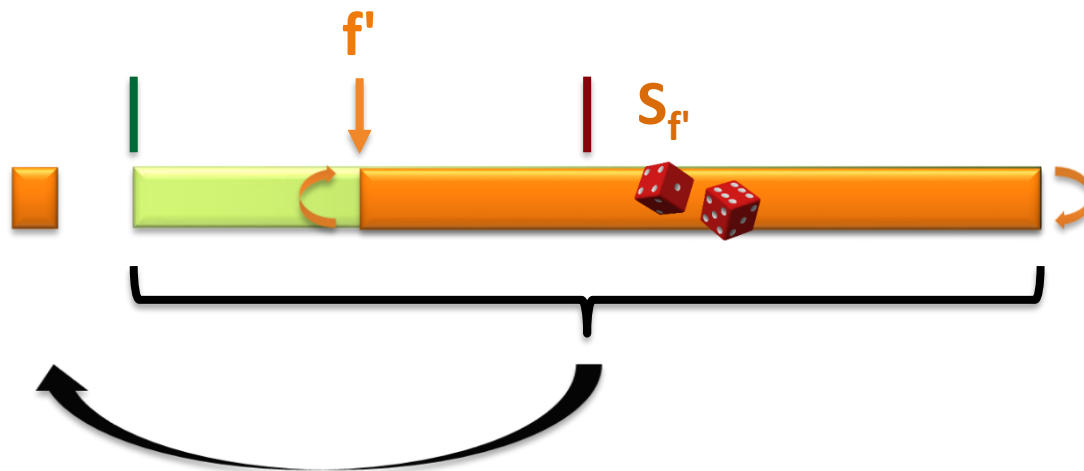


1.3 Find the largest subset of non-significant features

b) Find the feature f of lowest rank such that the subset of all features of higher ranks is not significant:

iii-a) accuracy \rightarrow or \nearrow

=> evaluate the performance after permutation of the features in $S_{f'}$



1.3 Find the largest subset of non-significant features

b) Find the feature f of lowest rank such that the subset of all features of higher ranks is not significant:

iii-b) accuracy \searrow

$\Rightarrow S_f$ contains significant features

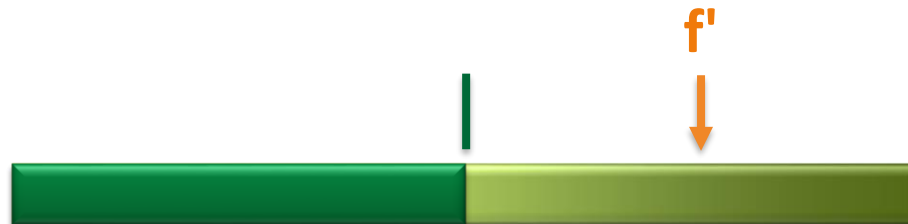


1.3 Find the largest subset of non-significant features

b) Find the feature f of lowest rank such that the subset of all features of higher ranks is not significant:

iii-b) accuracy \searrow

=> shift f downward to the mean rank of non-significant features

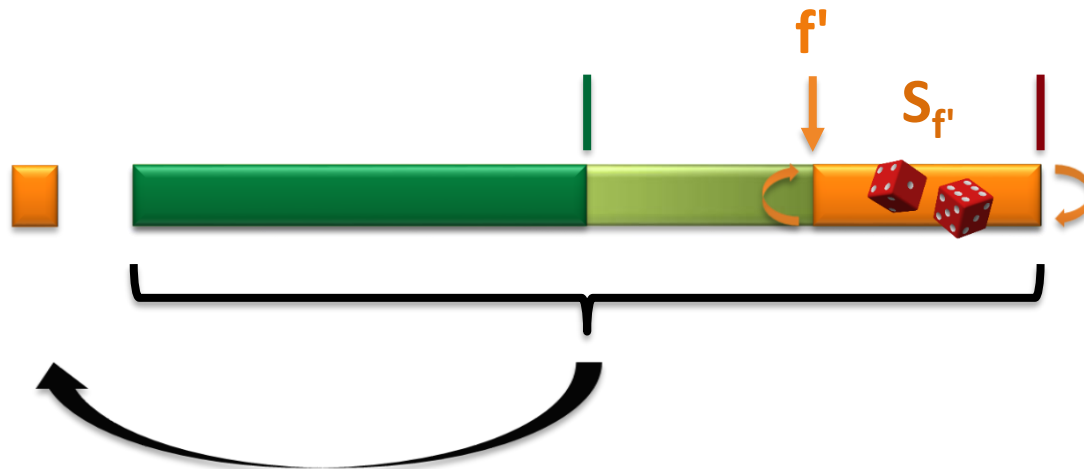


1.3 Find the largest subset of non-significant features

b) Find the feature f of lowest rank such that the subset of all features of higher ranks is not significant:

iii-b) accuracy \searrow

=> evaluate the performance after permutation of the features in $S_{f'}$



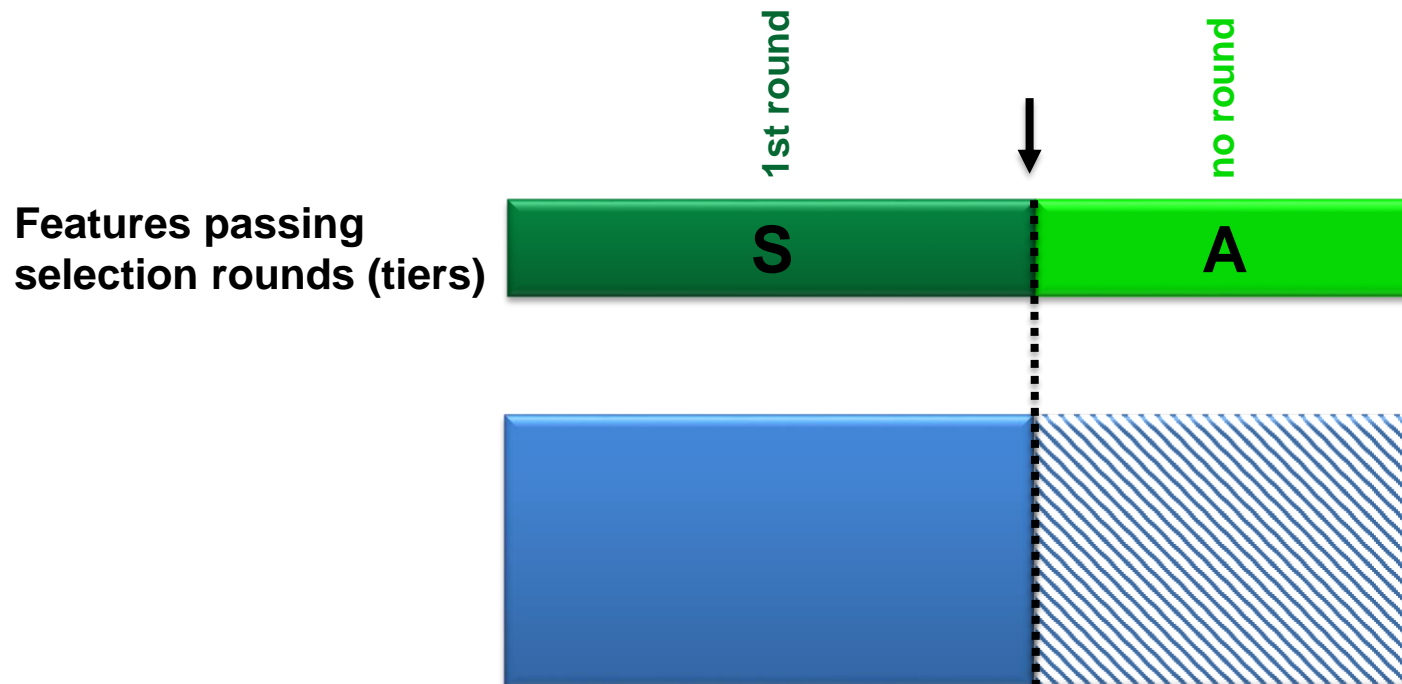
1.3 Find the largest subset of non-significant features

b) Find the feature f of lowest rank such that the subset of all features of higher ranks is not significant:

iv) stop when the upper and lower limits for f converge



1.4 Restrict the dataset to the significant subset

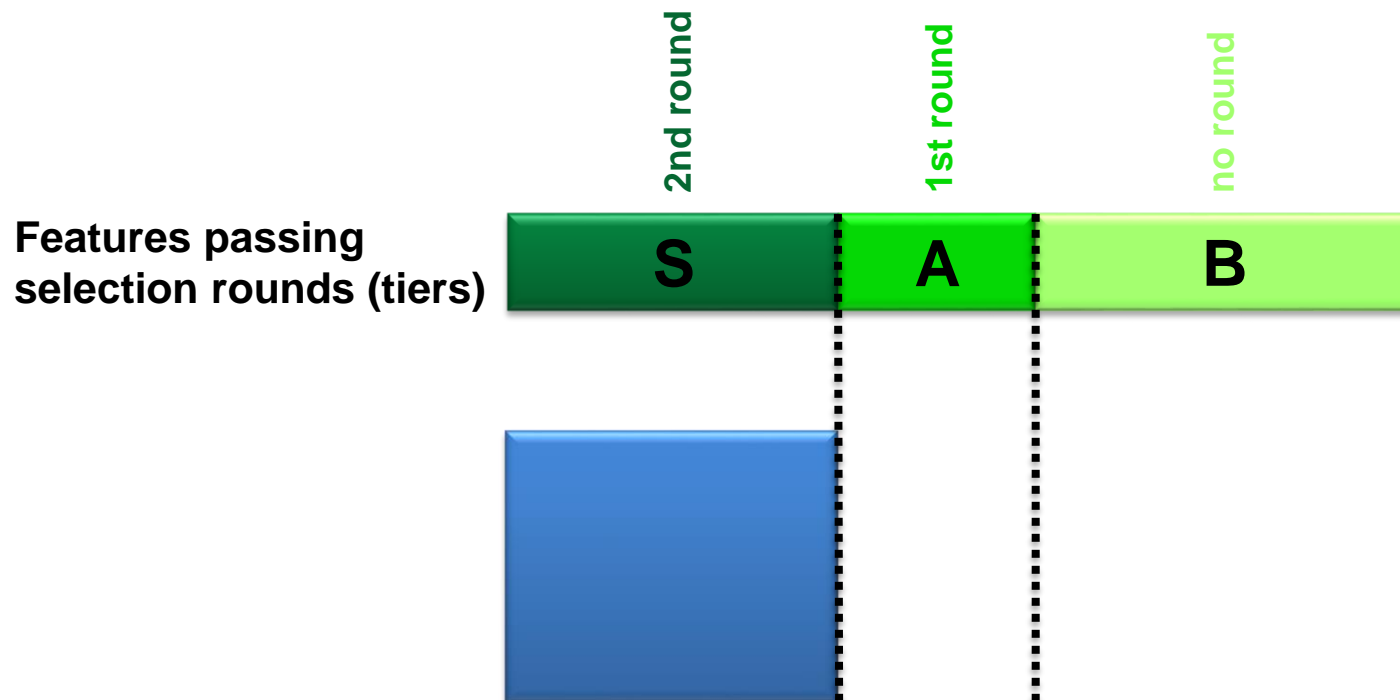


2. Repeat whole feature selection procedure on the restricted dataset

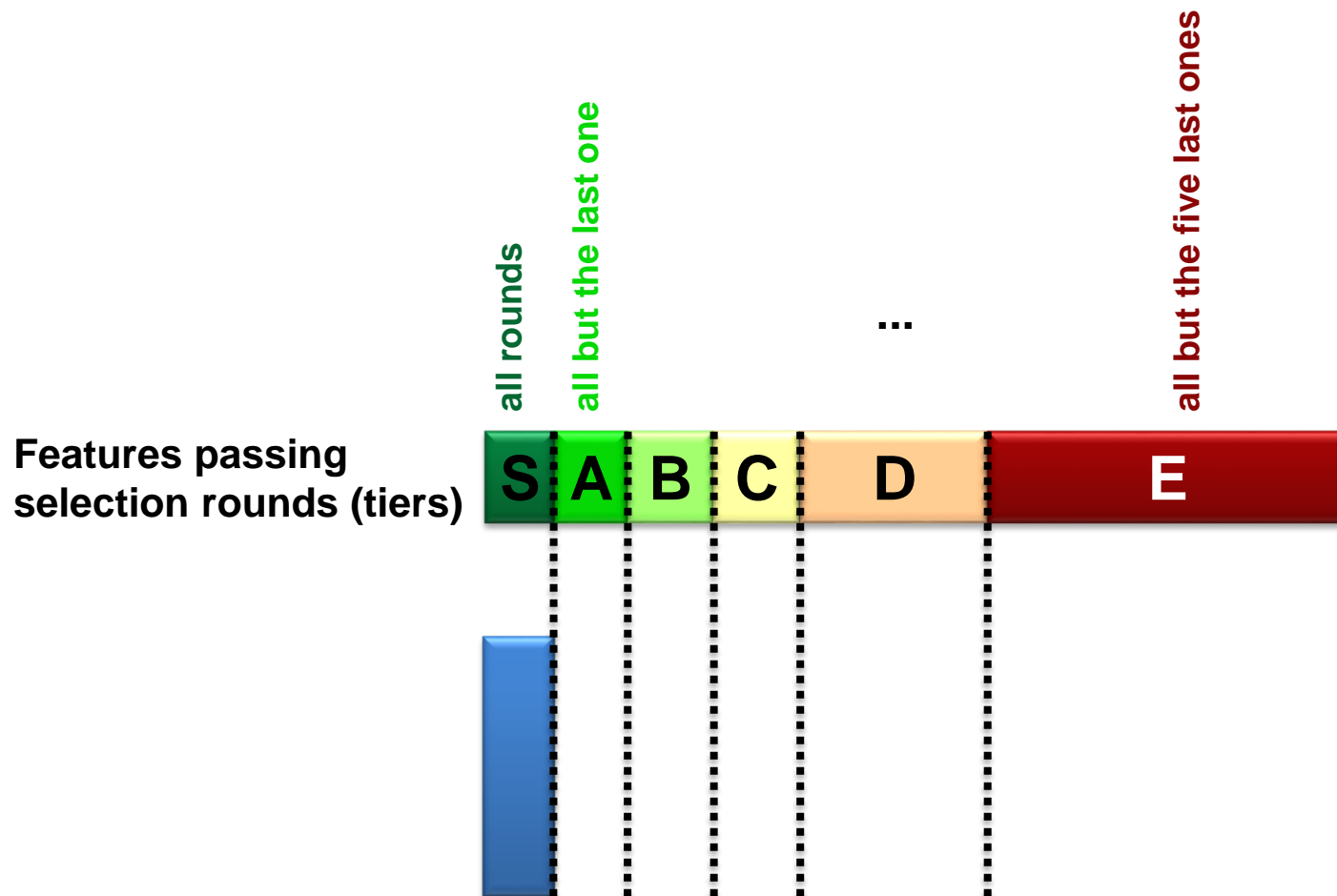
response (y) dataset(X')



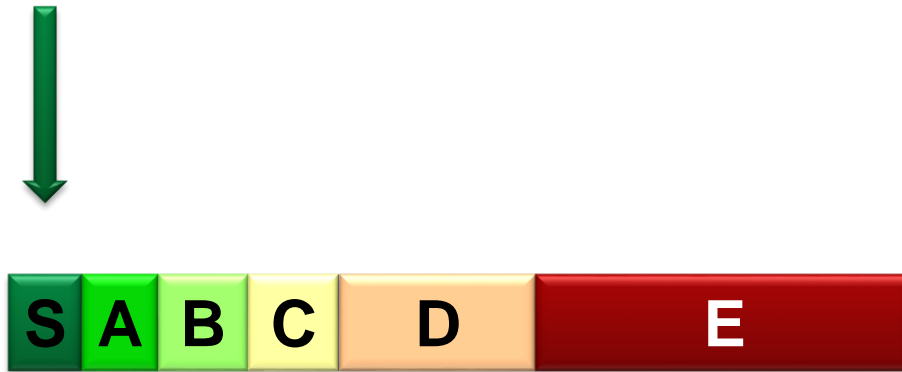
2. Repeat whole feature selection procedure on the restricted dataset



3. Stop when the signature is stable (all features significant)



Selected molecular signature



Final model



$$y = F_{\text{final}}(X_{\text{final}})$$

cea The biosigner software

➤ R package: Bioconductor
([DOI:10.18129/B9.bioc.biosigner](https://doi.org/10.18129/B9.bioc.biosigner))



➤ Galaxy tool: Toolshed, Workflow4Metabolomics, PhenoMeNal



➤ Publication: Frontiers in Molecular Biosciences
([DOI:10.3389/fmolb.2016.00026](https://doi.org/10.3389/fmolb.2016.00026))



ORIGINAL RESEARCH
published: 21 June 2016
doi: 10.3389/fmolb.2016.00026



***biosigner*: A New Method for the Discovery of Significant Molecular Signatures from Omics Data**

Philippe Rinaudo¹, Samia Boudah², Christophe Junot² and Etienne A. Thévenot^{1*}

Tools

Format Conversion

Preprocessing

Normalisation

Quality Control

Statistical Analysis

Anova N-way anova. With ou Without interactions

Hierarchical Clustering using ctc R package for java-treeview

Univariate Univariate statistics

Heatmap Heatmap of the dataMatrix

ACP ellipsoid by factors

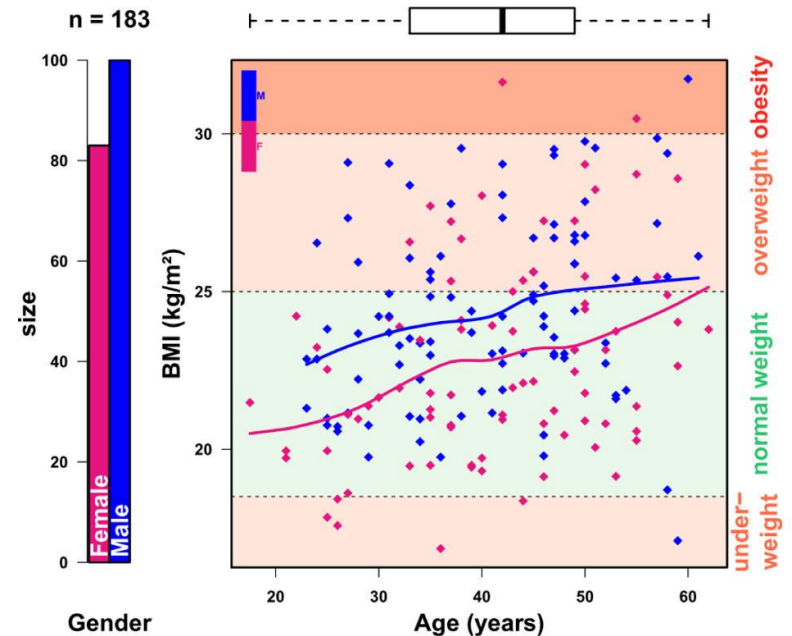
Biosigner Molecular signature discovery from omics data

Multivariate PCA, PLS and OPLS

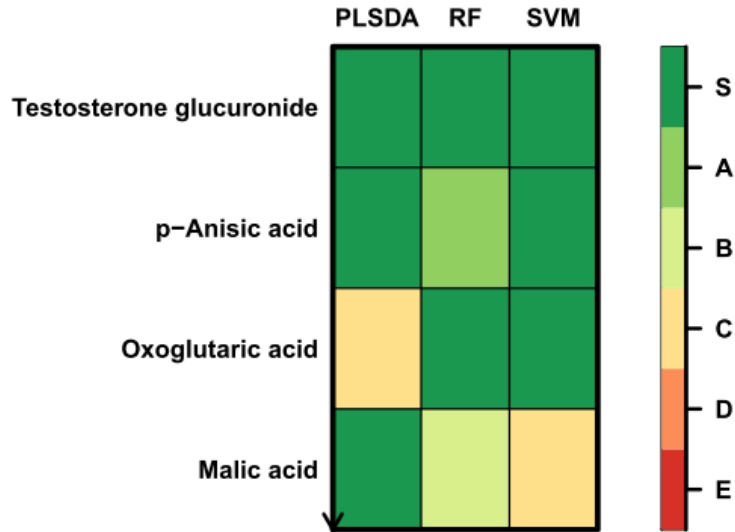
- Objective: influence of age, body mass index and gender on metabolite concentrations in urine
- Cohort: 184 employees from the CEA institute
- Analytics: LTQ-Orbitrap (negative ionization mode)
- Annotation: 109 metabolites were identified or annotated at the MSI level 1 or 2

- Pre-processing:

- XCMS followed by Quan Browser
- Signal drift and batch effect correction
- Normalization to the osmolality
- log10 transformation

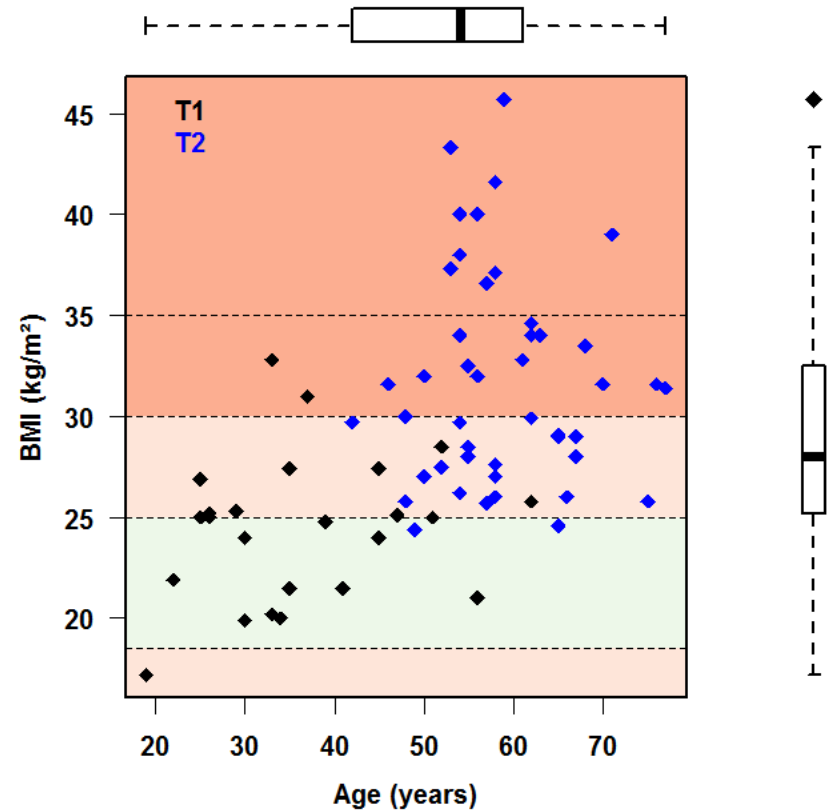


sacurine (ropIs)

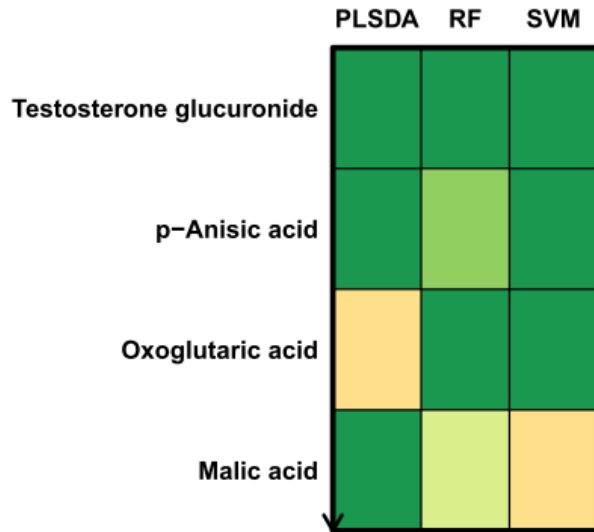


		sacurine
factor		gender
samples		183
features		109
signatures		[2-3]
performances (full -> restricted)	PLS-DA	87% -> 89%
	Random Forest	86% -> 86%
	SVM	88% -> 89%

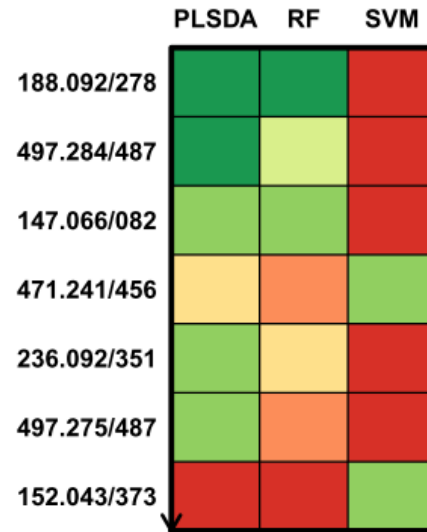
- LC-HRMS analysis of plasma
- from a cohort of 69 diabetic patients
- type 1 and type 2 patients
- 5,501 mz/RT features



sacurine (ropls)

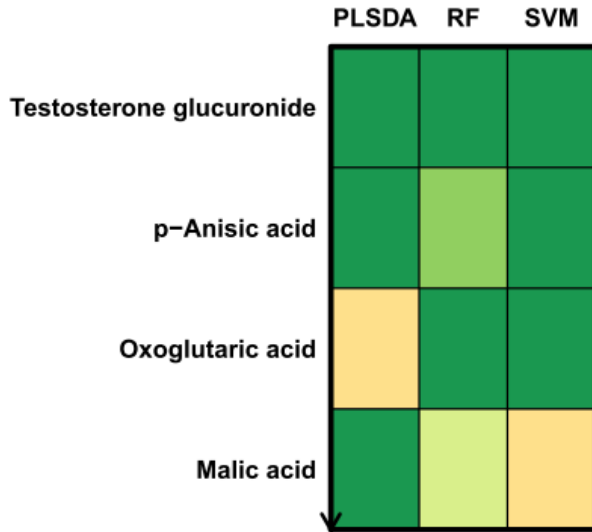


diaplasma (biosigner)

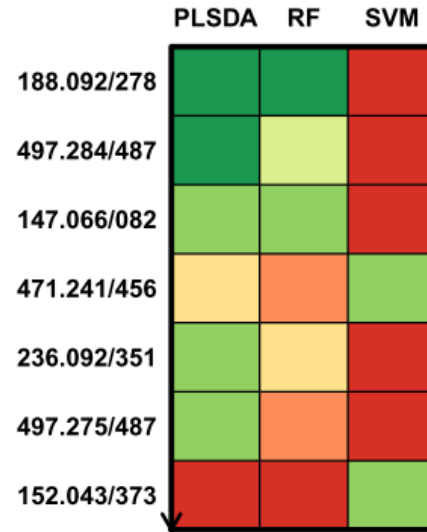


		sacurine	diaplasma
factor		gender	diabetic type
samples		183	69
features		109	5,501
signatures		[2-3]	[0-2]
performances (full -> restricted)	PLS-DA	87% -> 89%	83% -> 91%
	Random Forest	86% -> 86%	81% -> 81%
	SVM	88% -> 89%	83% -> na

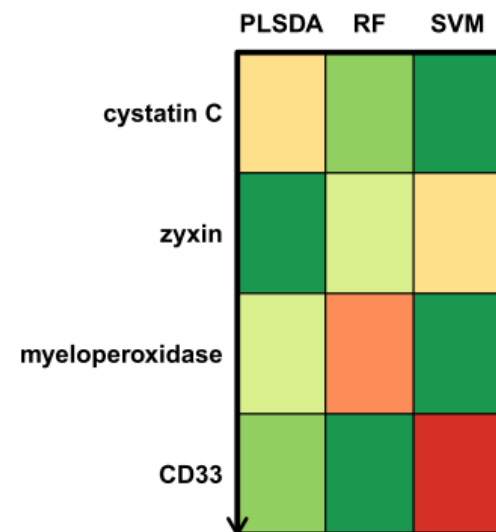
sacurine (ropls)



diaplasma (biosigner)



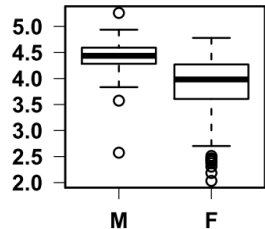
leukemia (golubEsets)



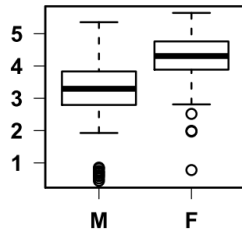
		sacurine	diaplasma	leukemia
factor		gender	diabetic type	ALL/AML
samples		183	69	72
features		109	5,501	7,129
signatures		[2-3]	[0-2]	[1-2]
performances (full -> restricted)	PLS-DA	87% -> 89%	83% -> 91%	95% -> 87%
	Random Forest	86% -> 86%	81% -> 81%	92% -> 92%
	SVM	88% -> 89%	83% -> na	93% -> 95%

sacurine

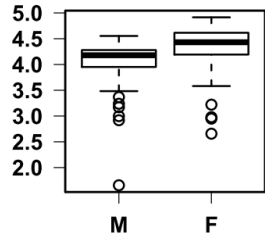
Testosterone glucuronide



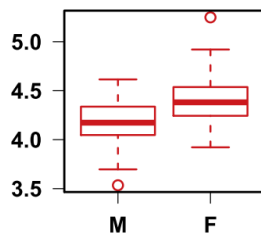
p-Anisic acid



Oxoglutaric acid

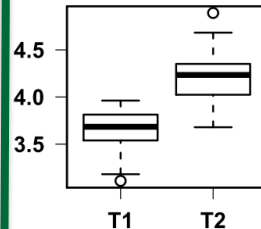


Malic acid

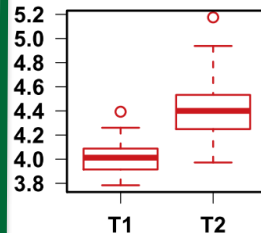


diaplasma

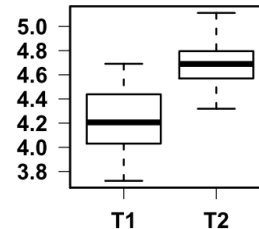
188.092/278



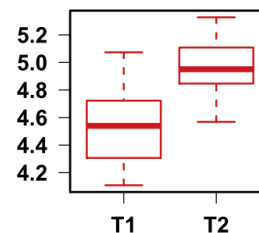
236.092/351



497.284/487

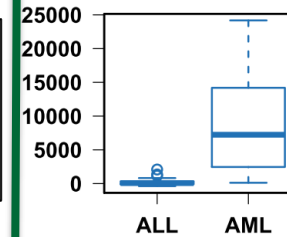


497.275/487

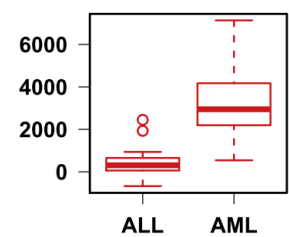


leukemia

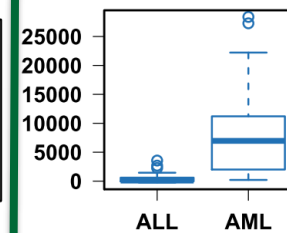
cystatin C



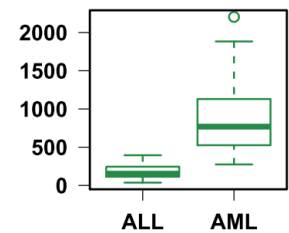
zyxin



myeloperoxidase



CD33

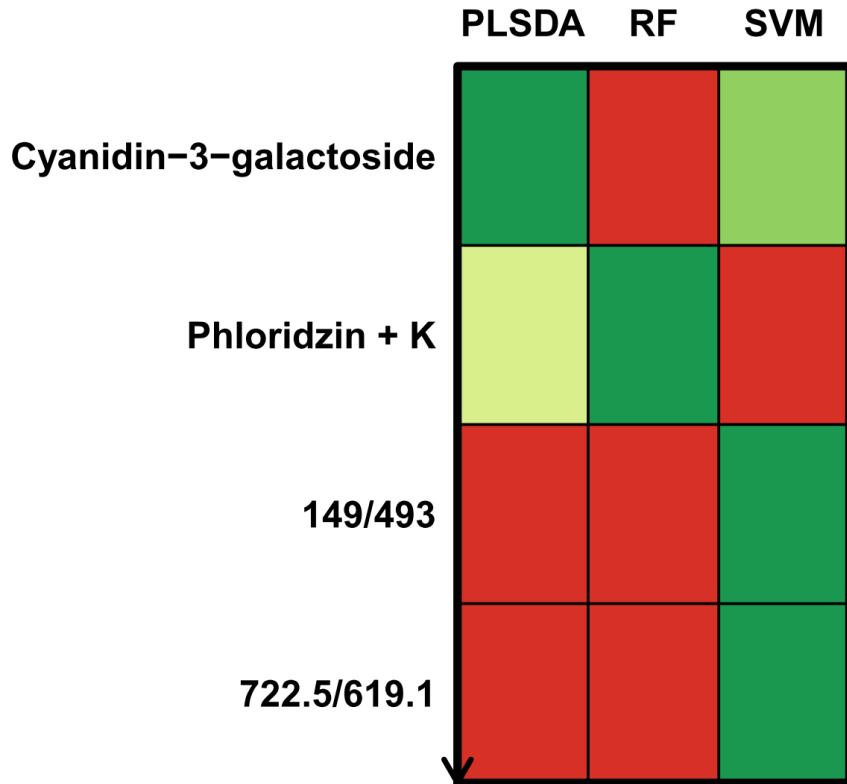


Biomarker in prostate cancer:
Zhang et al. (2013).
PLoS ONE, **8**:e65880.

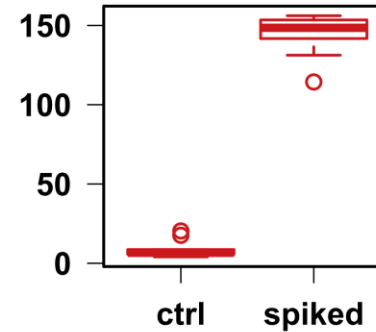
Taurochenodeoxycholic acid:
variation in type 2 diabetic patients:
Taylor et al. (2014). *PLoS ONE*,
9:e93540.

Cytochemical marker for
the diagnosis of AML:
Matsuo et al (2003).
Leukemia **17**:1538-1543.

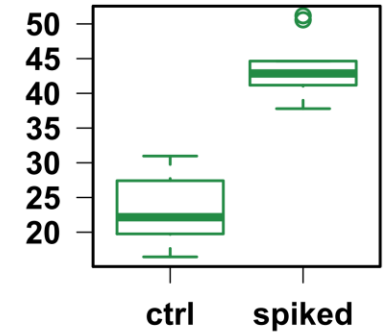
(BioMark)



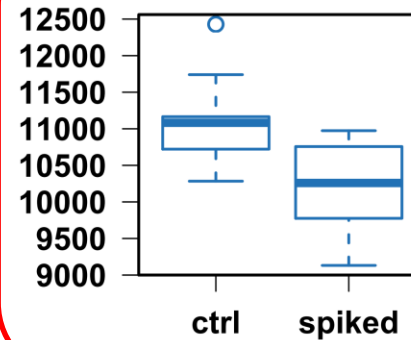
Cyanidin-3-galactoside



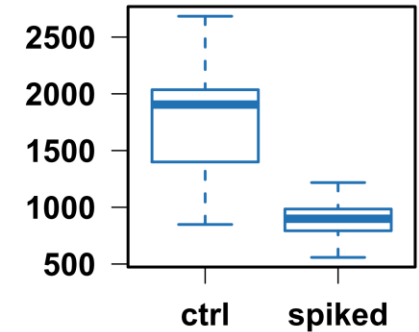
Phloridzin + K



149/493

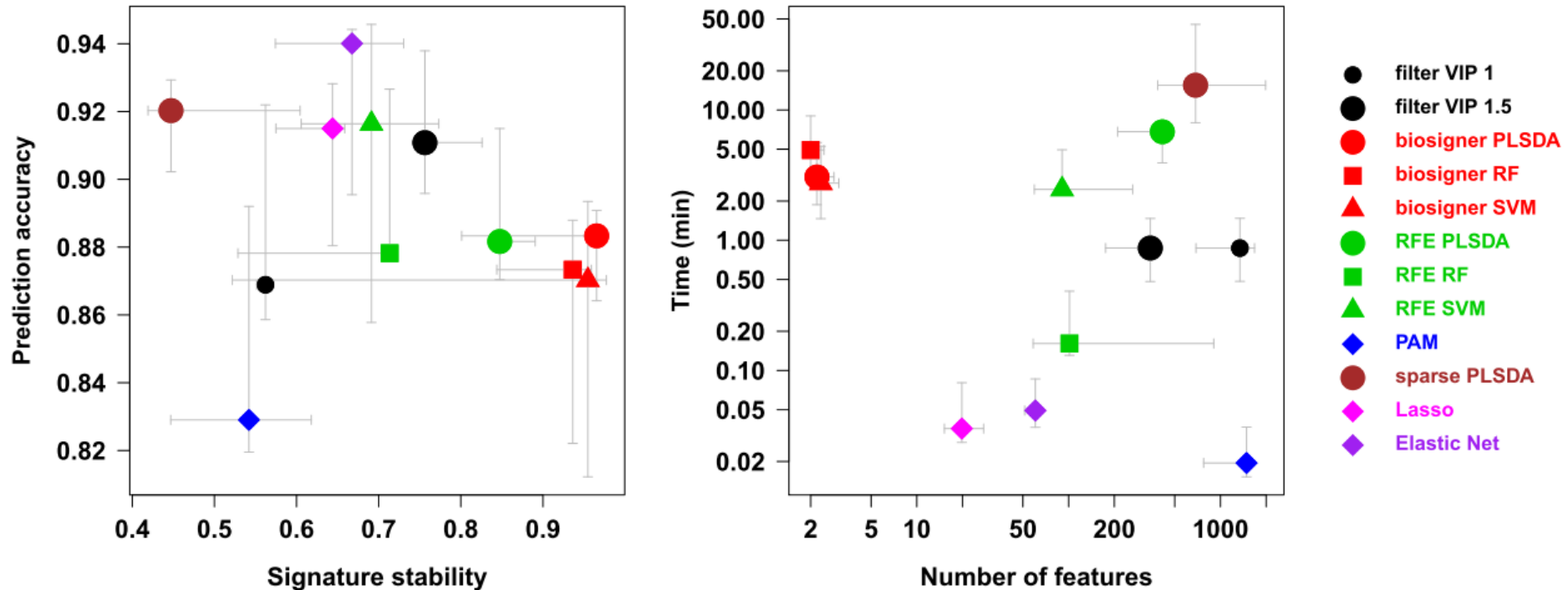


722.5/619.1



➤ SVM highlights features with decreased concentrations in spiked samples

Comparison with alternative feature selection methods



➤ *biosigner* finds small signatures providing a good compromise between prediction accuracy, signature stability and computation time

➤ Biosigner:

- efficient selection of significant signatures for binary classification
 - easy access (R and Galaxy)
-
- Depends on the structure of the dataset (distribution, limit of detection, correlation)
 - Validation on an independant dataset is mandatory
 - Importance of public datasets and code to benchmark new algorithms

➤ **Statistics**

- **Data integration**

ProMetIS: integration of proteomics and metabolomics data

➤ New methods and bioinformatics tools

- statistical integration (multivariate-based approaches)
- network analysis (pathway-based approaches)

➤ 2 case studies:

- high-throughput phenotyping of mouse models
- systems microbiology

➤ Large consortium

- CEA (LIST, IG, BIG), INRA (PFEM, TOXALIM), CNRS (LABGeM)

➤ 2 year project (« Integrative Bioinformatics » workpackages from IFB)

➤ Perspectives:

- application to human phenotyping (France Medecine Genomique 2025)

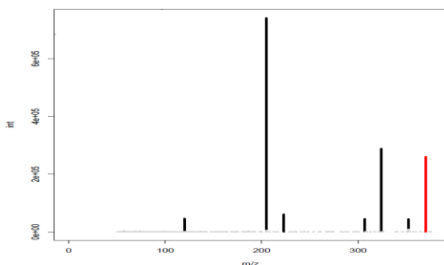




cea Spectrum identification

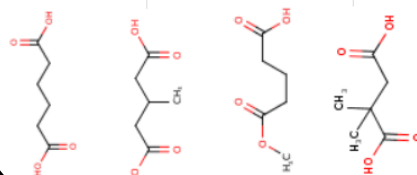
Input:

- A mass spectrum
- A precursor m/z



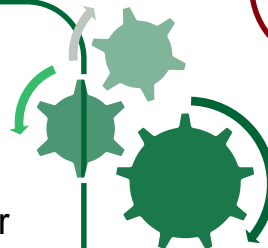
- CSI-FingerID (Shen, 2014)
- CFM-ID (Allen, 2014)
- MetFrag (Wolf, 2010)
- Etc...

Molecular and/or spectral **database**

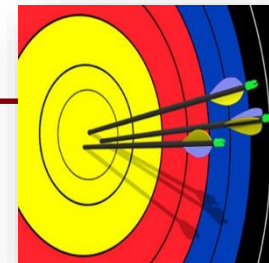
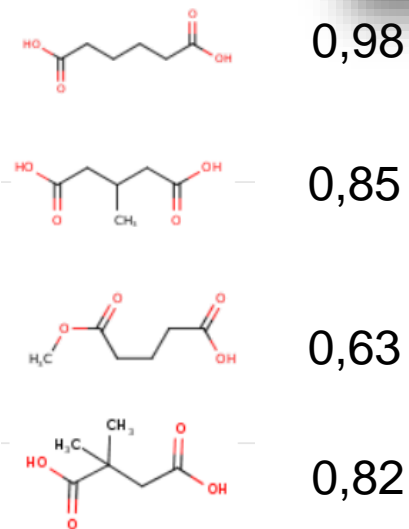


Identification software*

Model learned on a database.
CFM-ID: Transition probabilities
CSI-FingerID: Linear transformation between kernel spaces



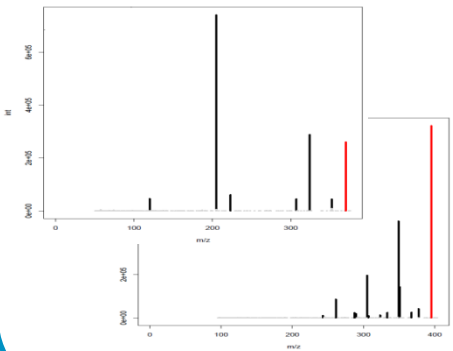
Output :
Scored molecules



cea Mining spectral libraries

Input:

- A set of spectra

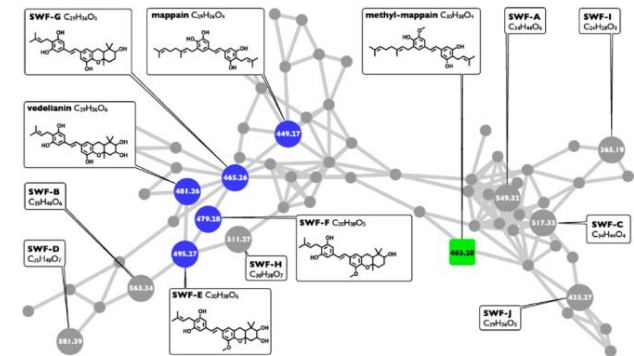


Spectral Mining software

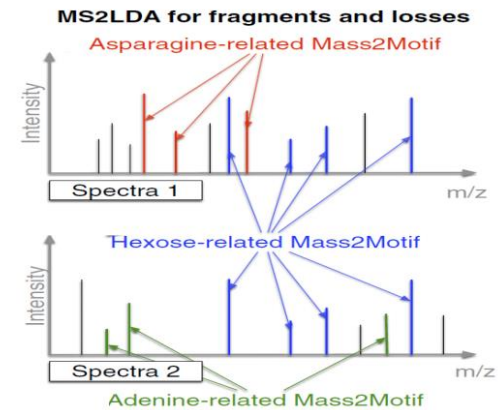
Output:

Information about structural similarities

- **Networks** (GNPS, Wang, 2014)



- **Motifs discovery** (MS2LDA, Van der Hooft, 2016)

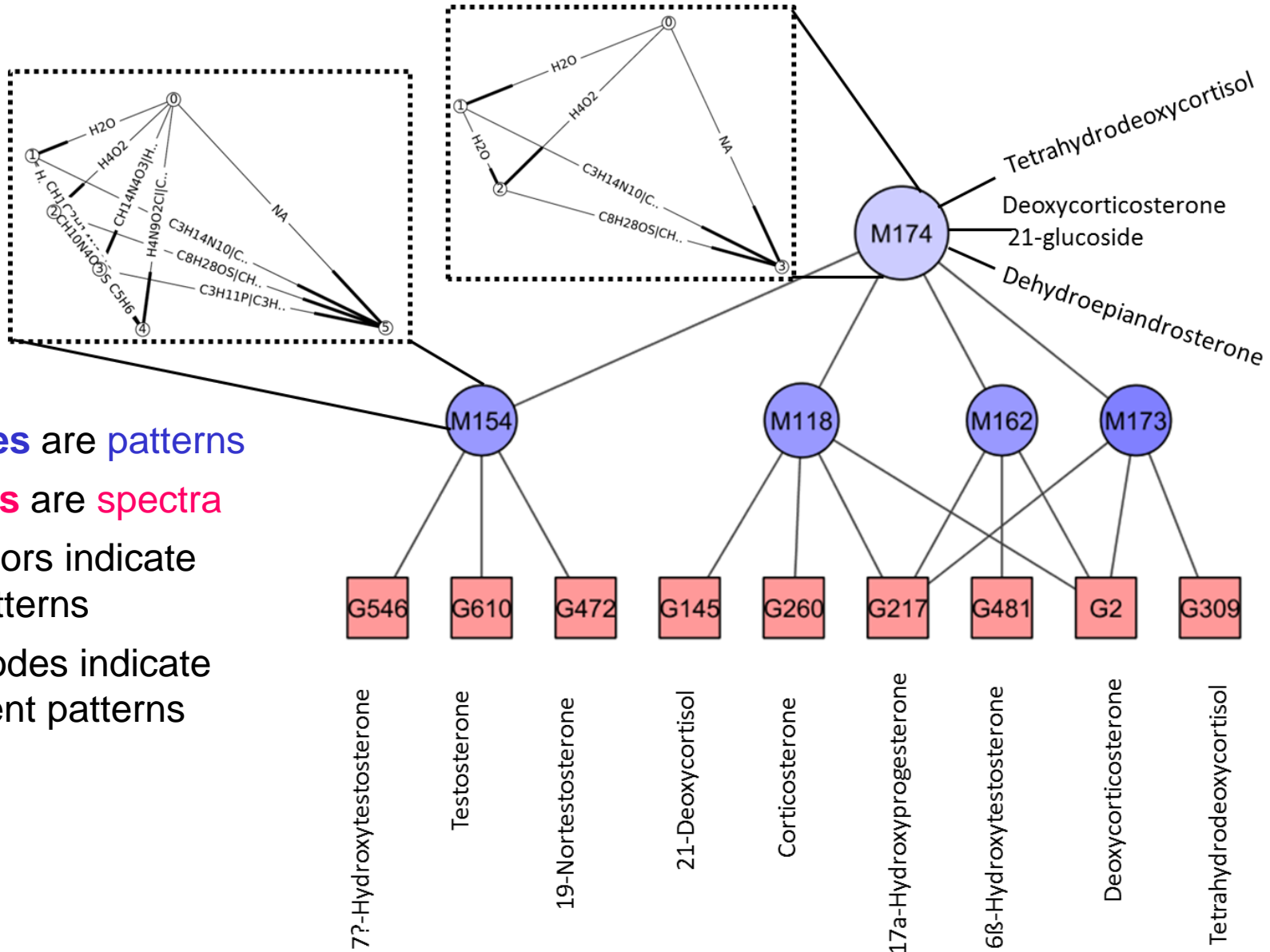


Running time and number of patterns

Dataset	<i>Penicillium verrucosum</i>	Reference library from pure compounds
Reference	Hautbergue <i>et al.</i> 2017 <i>J Chromatogr B</i>	Metabolome IDF
Number of spectra	91	834
Graph building	35 s	1 min 45 s
Pattern Mining	10 s	2 min 10 s
Number of patterns	54	832

- Processing time is more dependent on the similarities in the datasets than of the size in the dataset.

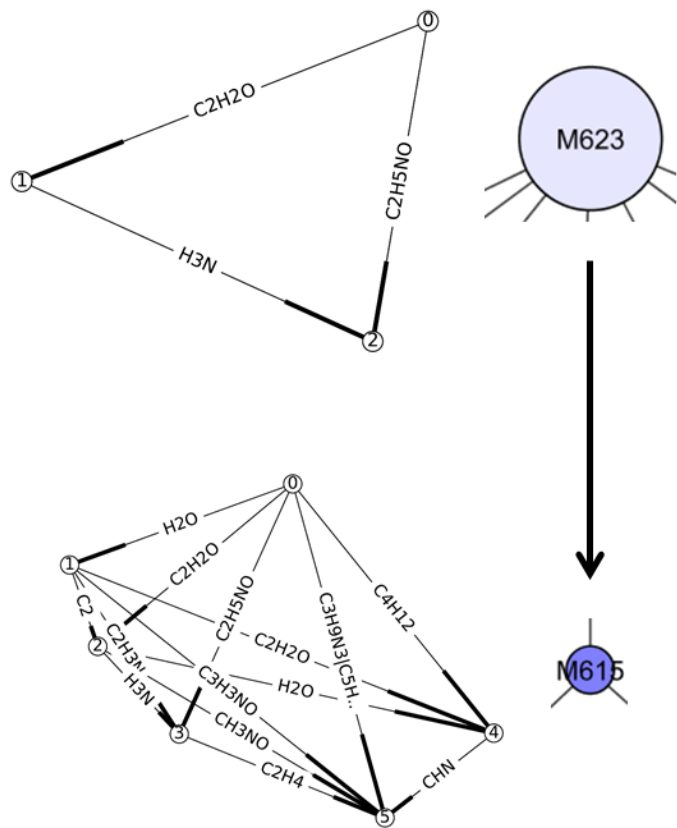
Example: sterone sub-lattice



- Blue nodes are patterns
- Red nodes are spectra
- Lighter colors indicate smaller patterns
- Smaller nodes indicate less frequent patterns

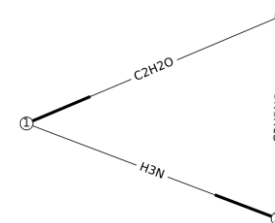
cea Example: sterone sub-lattice

Patterns include coarse and fine grain information

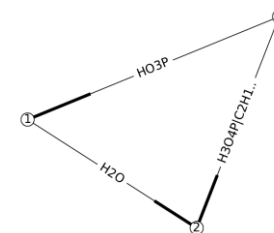


➤ Example of generic patterns:

■ N-Acetyl * compounds



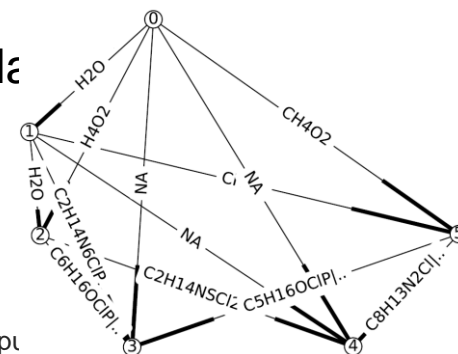
■ Phosphated compounds



■

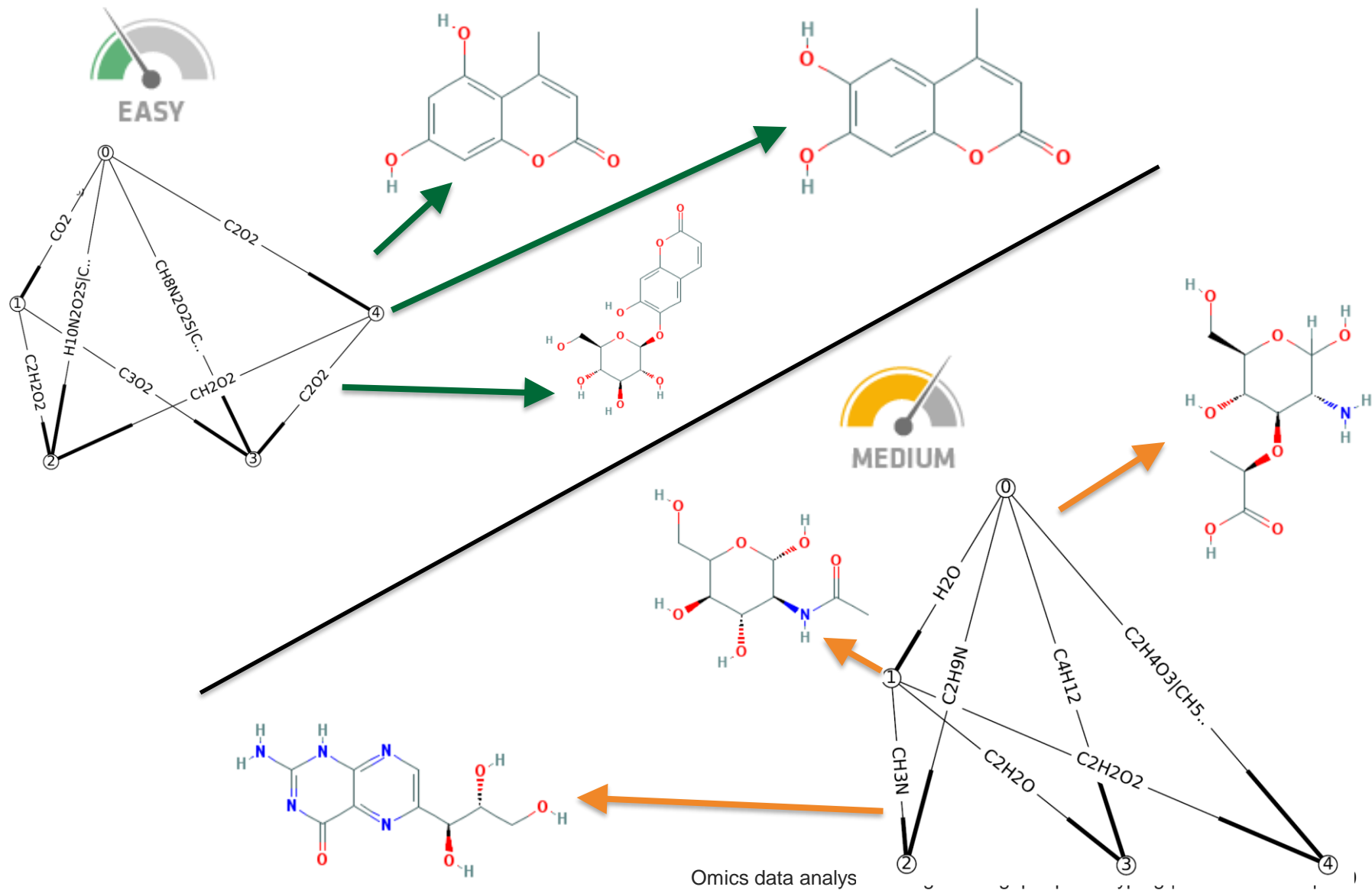
➤ Specific patterns:

■ Sphingosine rela



Major difference with the patterns obtained with Mass2LDA (Van der Hooft *et al.* 2016, *PNAS*)

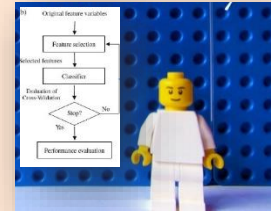
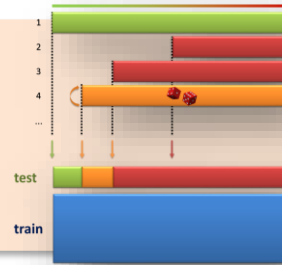
cea Relating patterns to chemical (sub-)structures



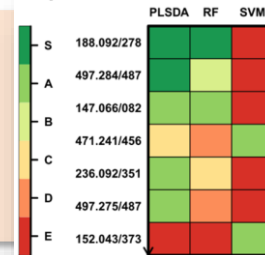
- **Workflow management**
 - **Workflows**

Implementation: From the method to the workflow

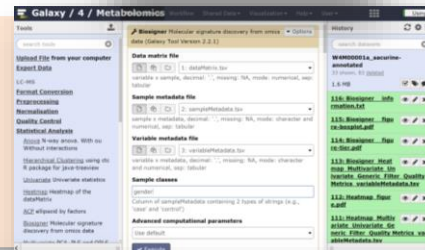
Method



Package

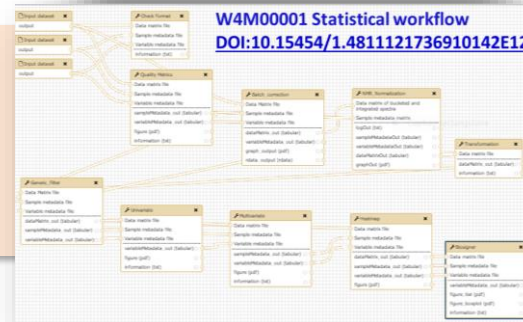


GUI



Workflow

Cloud



<http://workflow4metabolomics.org>

The workflow challenge: bridging experimenters' and bioinformaticians' talents

Users

- **Web-based**
- **Workflow editor**
- **User-friendly**
- **Tutorials**



- **Shared workflows**
- **Reproducible science**



Developers

- **Multi-language**
- **Toolshed**
- **Open-source**



- **Workflow management**
 - **Galaxy environment**

<https://galaxyproject.org>

- Workflow management through a classic web browser

[Giardine *et al.* \(2005\). Galaxy: A platform for interactive large-scale genome analysis. *Genome Res.* 15:1451-1455.](#)

- Started in 2005; more than 55,000 users worldwide

[Goecks *et al.* \(2010\). Galaxy: a comprehensive approach for supporting accessible, reproducible, and transparent computational research in the life sciences. *Genome Biol.* 11:R86.](#)

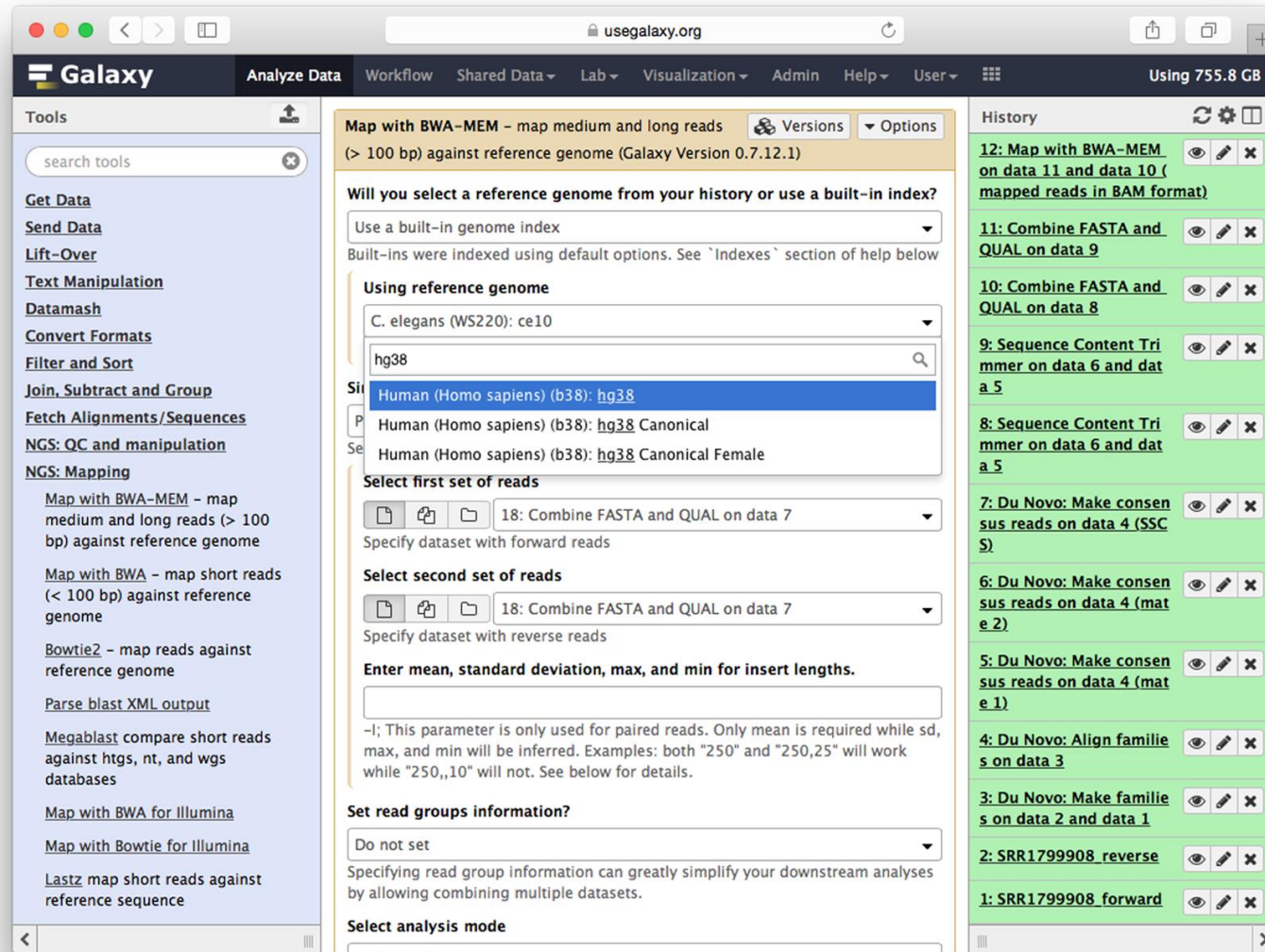
[Afgan *et al.* \(2016\). The Galaxy platform for accessible, reproducible and collaborative biomedical analyses: 2016 update. *Nature Biotechnol.* 29:972-974.](#)

- NGS, transcriptomics, proteomics, metabolomics

[Boekel *et al.* \(2015\). Multi-omic data analysis using Galaxy. *Nature Biotechnol.* 33:137-139.](#)

[Guitton *et al.* \(2017\). Create, run, share, publish, and reference your LC-MS, FIA-MS, GC-MS, and NMR data analysis workflows with the Workflow4Metabolomics 3.0 Galaxy online infrastructure for metabolomics. *Int. J. Biochem. Cell. Biol.* 93:89-101.](#)





The screenshot displays the Galaxy web interface for configuring the 'Map with BWA-MEM' tool. The main panel shows the tool title and version information, followed by a question about selecting a reference genome. A dropdown menu is open, showing 'hg38' as the selected option. Below this, there are sections for selecting the first and second sets of reads, and a section for entering mean, standard deviation, max, and min for insert lengths. The right sidebar shows a history of previous tool runs, with the current run highlighted in green.

Galaxy Analyze Data Workflow Shared Data Lab Visualization Admin Help User Using 755.8 GB

Tools

search tools

Get Data

Send Data

Lift-Over

Text Manipulation

Datamash

Convert Formats

Filter and Sort

Join, Subtract and Group

Fetch Alignments/Sequences

NGS: QC and manipulation

NGS: Mapping

Map with BWA-MEM - map medium and long reads (> 100 bp) against reference genome

Map with BWA - map short reads (< 100 bp) against reference genome

Bowtie2 - map reads against reference genome

Parse blast XML output

Megablast compare short reads against htgs, nt, and wgs databases

Map with BWA for Illumina

Map with Bowtie for Illumina

Lastz map short reads against reference sequence

Map with BWA-MEM - map medium and long reads Versions Options

(> 100 bp) against reference genome (Galaxy Version 0.7.12.1)

Will you select a reference genome from your history or use a built-in index?

Use a built-in genome index

Built-ins were indexed using default options. See `Indexes` section of help below

Using reference genome

C. elegans (WS220): ce10

hg38

Human (Homo sapiens) (b38): **hg38**

Human (Homo sapiens) (b38): hg38 Canonical

Human (Homo sapiens) (b38): hg38 Canonical Female

Select first set of reads

18: Combine FASTA and QUAL on data 7

Specify dataset with forward reads

Select second set of reads

18: Combine FASTA and QUAL on data 7

Specify dataset with reverse reads

Enter mean, standard deviation, max, and min for insert lengths.

-l; This parameter is only used for paired reads. Only mean is required while sd, max, and min will be inferred. Examples: both "250" and "250,25" will work while "250,,10" will not. See below for details.

Set read groups information?

Do not set

Specifying read group information can greatly simplify your downstream analyses by allowing combining multiple datasets.

Select analysis mode

History

12: **Map with BWA-MEM on data 11 and data 10 (mapped reads in BAM format)**

11: **Combine FASTA and QUAL on data 9**

10: **Combine FASTA and QUAL on data 8**

9: **Sequence Content Trimmer on data 6 and data 5**

8: **Sequence Content Trimmer on data 6 and data 5**

7: **Du Novo: Make consensus reads on data 4 (SSC S)**

6: **Du Novo: Make consensus reads on data 4 (mate 2)**

5: **Du Novo: Make consensus reads on data 4 (mate 1)**

4: **Du Novo: Align families on data 3**

3: **Du Novo: Make families on data 2 and data 1**

2: **SRR1799908_reverse**

1: **SRR1799908_forward**

[Afgan et al. \(2016\). The Galaxy platform for accessible, reproducible and collaborative biomedical analyses: 2016 update. Nat. Biotechnol. 29:972-974.](#)

tools

Upload File from your computer

Export Data

LC-MS

Format Conversion

Preprocessing

Normalisation

Quality Control

Statistical Analysis

Annotation

GC-MS

Preprocessing

Normalisation

Quality Control

Statistical Analysis

Annotation

NMR

Preprocessing

Normalisation

Quality Control

Statistical Analysis

Annotation

COMMON TOOLS

Data Handling

Text Manipulation

Filter and Sort

Join, Subtract and Group

Statistics

Graph/Display Data

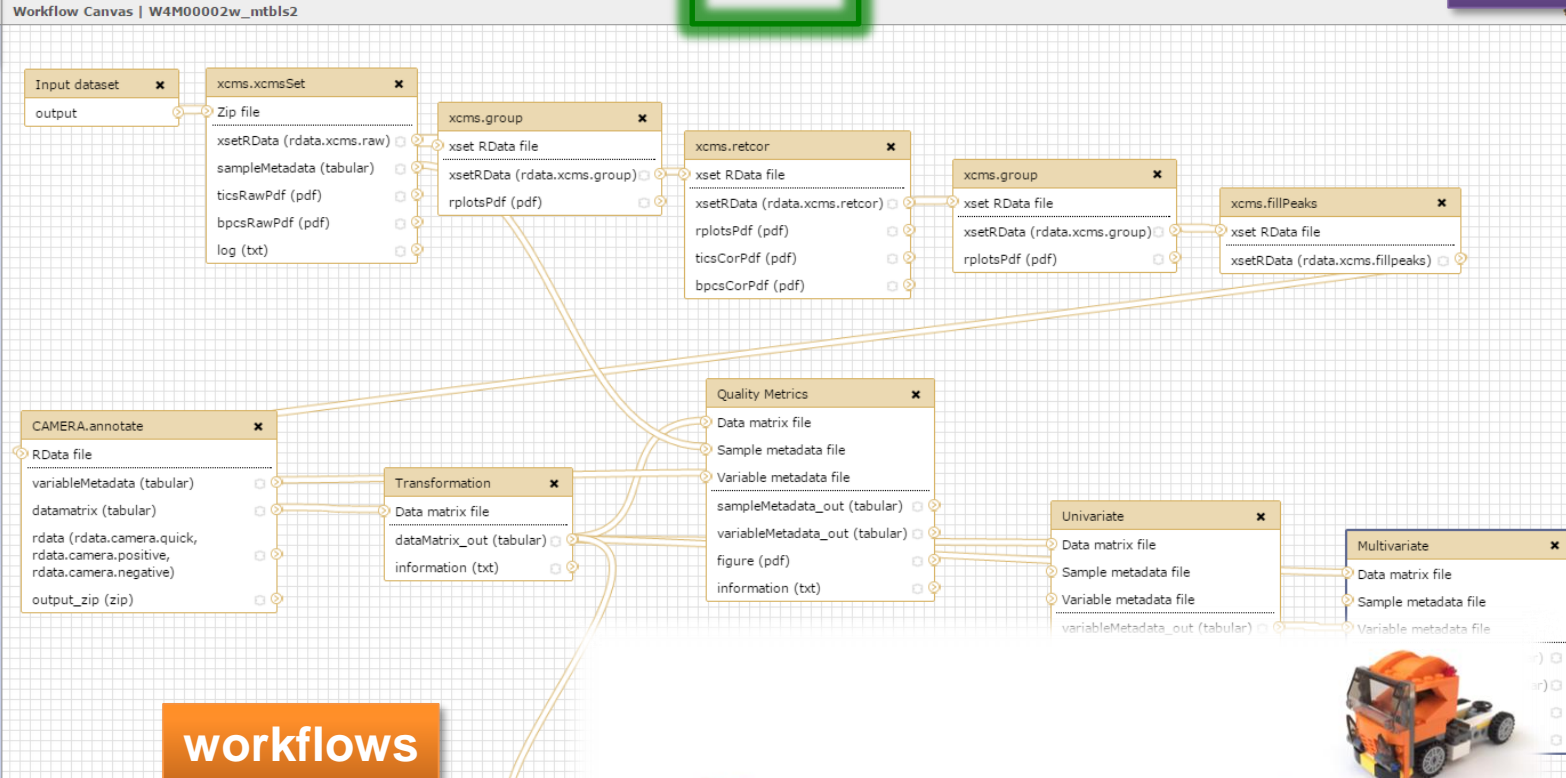
Deprecated Tools

New tools Version

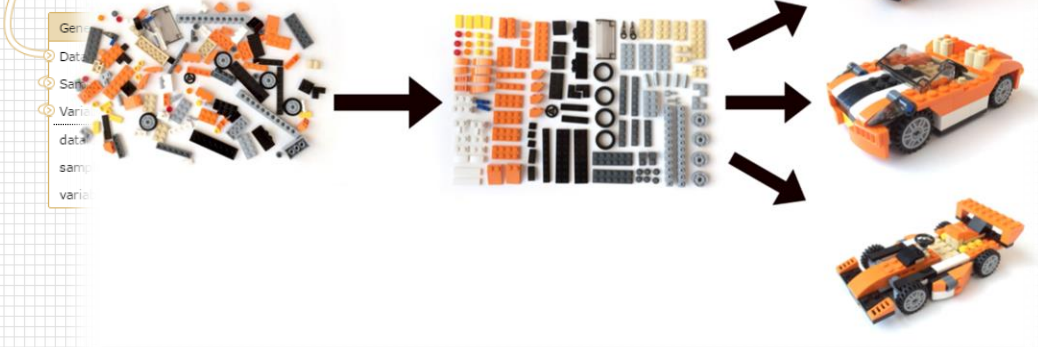
Multiple regression

Workflow control

Inputs



workflows



Details

Multivariate

PCA, PLS and OPLS (Galaxy Tool Version 2.2.4)

Data matrix file

Data input 'dataMatrix_in' (tabular) variable x sample, decimal: '.', missing: NA, mode: numerical, sep: tabular

Sample metadata file

Data input 'sampleMetadata_in' (tabular) sample x metadata, decimal: '.', missing: NA, mode: character and numerical, sep: tabular

Variable metadata file

Data input 'variableMetadata_in' (tabular) variable x metadata, decimal: '.', missing: NA, mode: character and numerical, sep: tabular

Y Response (for PLS(-DA) and OPLS(-DA) only)

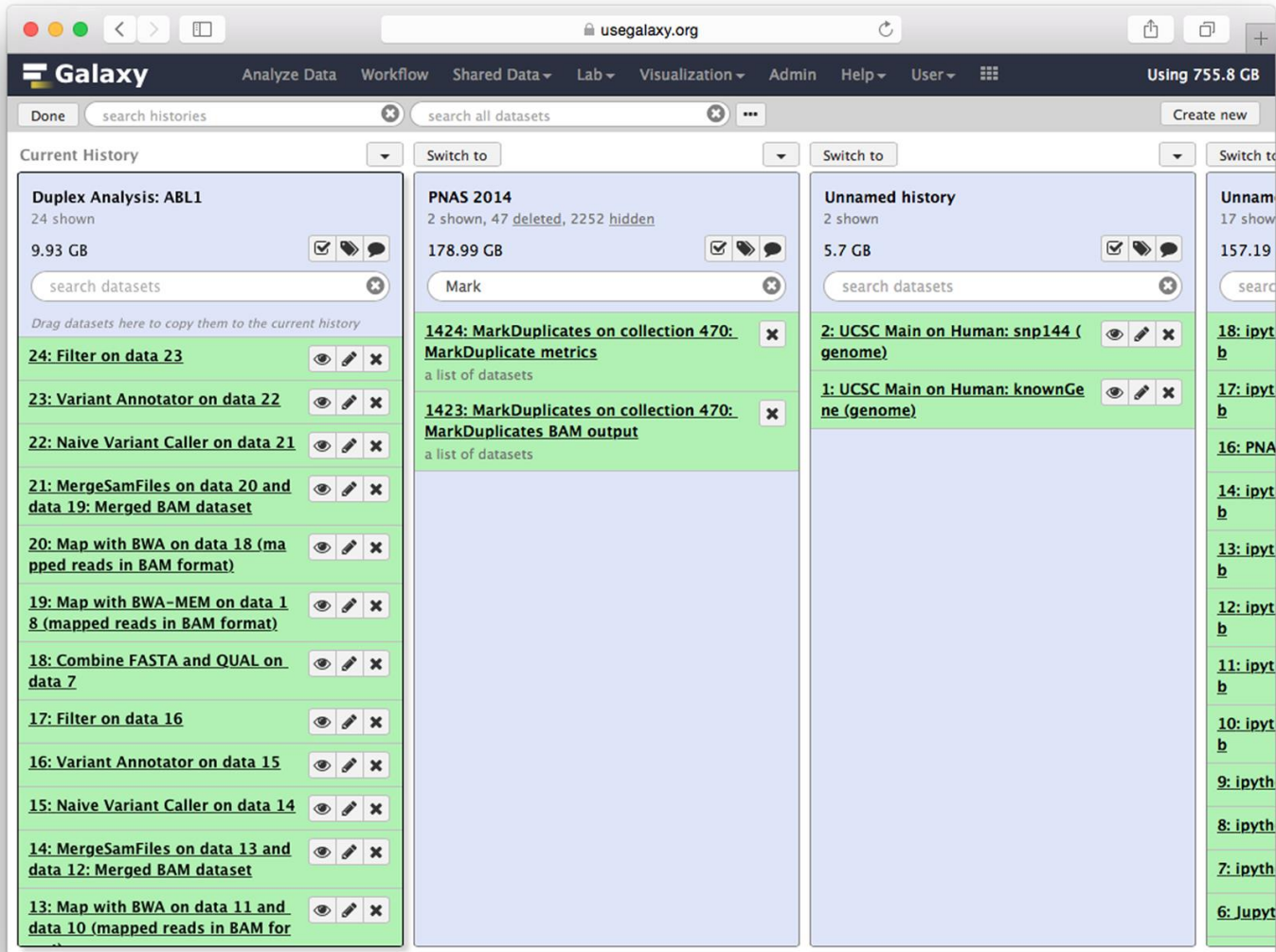
class

Notes: 1) PCA: keep the default (none); 2) PLS(-DA) and OPLS(-DA): indicate the name of the column of the sample table to be modeled

Number of predictive components

NA

Notes: 1) PCA and PLS(-DA): NA can be selected to get a



[Afgan *et al.* \(2016\). The Galaxy platform for accessible, reproducible and collaborative biomedical analyses: 2016 update. *Nat. Biotechnol.* 29:972-974.](#)

Share or Publish History 'W4M00001_Sacurine-statistics'

Make History Accessible via Link and Publish It

This history is currently **accessible via link and published**.
Anyone can view and import this history by visiting the following URL:

<http://galaxy.workflow4metabolomics.org/u/ethevenot/h/w4m00001sacurine-statistics-1>

This history is publicly listed and searchable in Galaxy's Published Histories section.
You can:

Unpublish History

Removes this history from Galaxy's Published Histories section so that it is not publicly listed or searchable.

Disable Access to History via Link and Unpublish

Disables this history's link so that it is not accessible and removes history from Galaxy's Published Histories section so that it is not publicly listed or searchable.

Share History with Individual Users

You have not shared this history with any users.

Share with a user

History ↻ ⚙️ 🗑️

search datasets ✖️

W4M00001_Sacurine-statistics
53 shown

3.99 MB ☑️ 📁 💬

53: Heatmap figure.pdf 👁️ ✎️ ✖️

52: Heatmap Multivariate Multivariate Univariate Univariate Generic Filter Quality Metrics Generic Filter Quality Metrics Batch correction all loess pool variable Metadata.tsv 👁️ ✎️ ✖️

51: Heatmap Multivariate Multivariate Univariate Generic Filter Quality Metrics Generic Filter Quality Metrics sampleMetadata.tsv 👁️ ✎️ ✖️

50: Heatmap Generic Filter Transform 👁️ ✎️ ✖️



- **Workflow management**
 - **Workflow4Metabolomics**
online platform

Workflow4metabolomics

Main menu

- Home
- Events
- History
- ▼ Introduction
 - The Galaxy environment
 - ▶ The LC-MS workflow
 - The GC-MS workflow
 - The NMR workflow
 - References
- HowTo
- ▼ Download
 - ▶ Datasets
- Referenced WorkFlows and Histories
- How to contribute?
- ▼ Developer resources
 - Source-code
 - Virtual environments

Workflow4Metabolomics 3.0

Welcome to the collaborative portal dedicated to metabolomics data processing, analysis and annotation for Metabolomics community.

" We are happy to announce the next **Workflow4Experimenters (W4E) international course 2018**: *Using Galaxy and the Workflow4metabolomics infrastructure to analyse metabolomics data.*
Please save the date: **8-12 October 2018** at Pasteur Institute, Paris - France

More news in April ! "

Follow us on Twitter  @workflow4metabo

STEP 1

STEP 2

STEP 3

STEP 4

[Giacomoni et al. \(2015\). Workflow4Metabolomics: a collaborative research infrastructure for computational metabolomics. *Bioinformatics*. 31:1493-1495.](#)

- **Coordinators (C. Médigue & J. van Helden)**



- French node from Elixir
- Federation of 34 national bioinformatics platforms
- 230 bioinformaticians

- **Missions**

- **E-infrastructure** components : Storage, Computing (e.g. **cloud**) , Tools, VRE...
- **Training** (NGS, Galaxy) & community animation
- Collaboration with national and European infrastructures



- **National task force**

- IFB Galaxy Working Group
- European Galaxy Developer workshop 2017
- Organization of the GCC conference 2017





METABOHUB

MetaboHUB: The French infrastructure for metabolomics and fluxomics

* Coordinator (D. Rolin)

* 80 permanent scientists

* Total budget: 45 M€

* Launched in 2013



* 4 LC-MS, GC-MS and NMR platforms

* dedicated to Innovation, Service, Technology Transfer and Training

* Built upon the Francophone Network for Metabolomics and Fluxomics (> 300 members)

* 4 online bioinformatics infrastructures

* workflows, databases, pathways

* Partnerships:



* bioinformatics, proteomics, cohorts, crops

* Europe:



* MetaboMOOC (oct. 2018; with IFB)  MOOC

Human Health



Agriculture



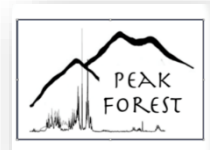
Environment



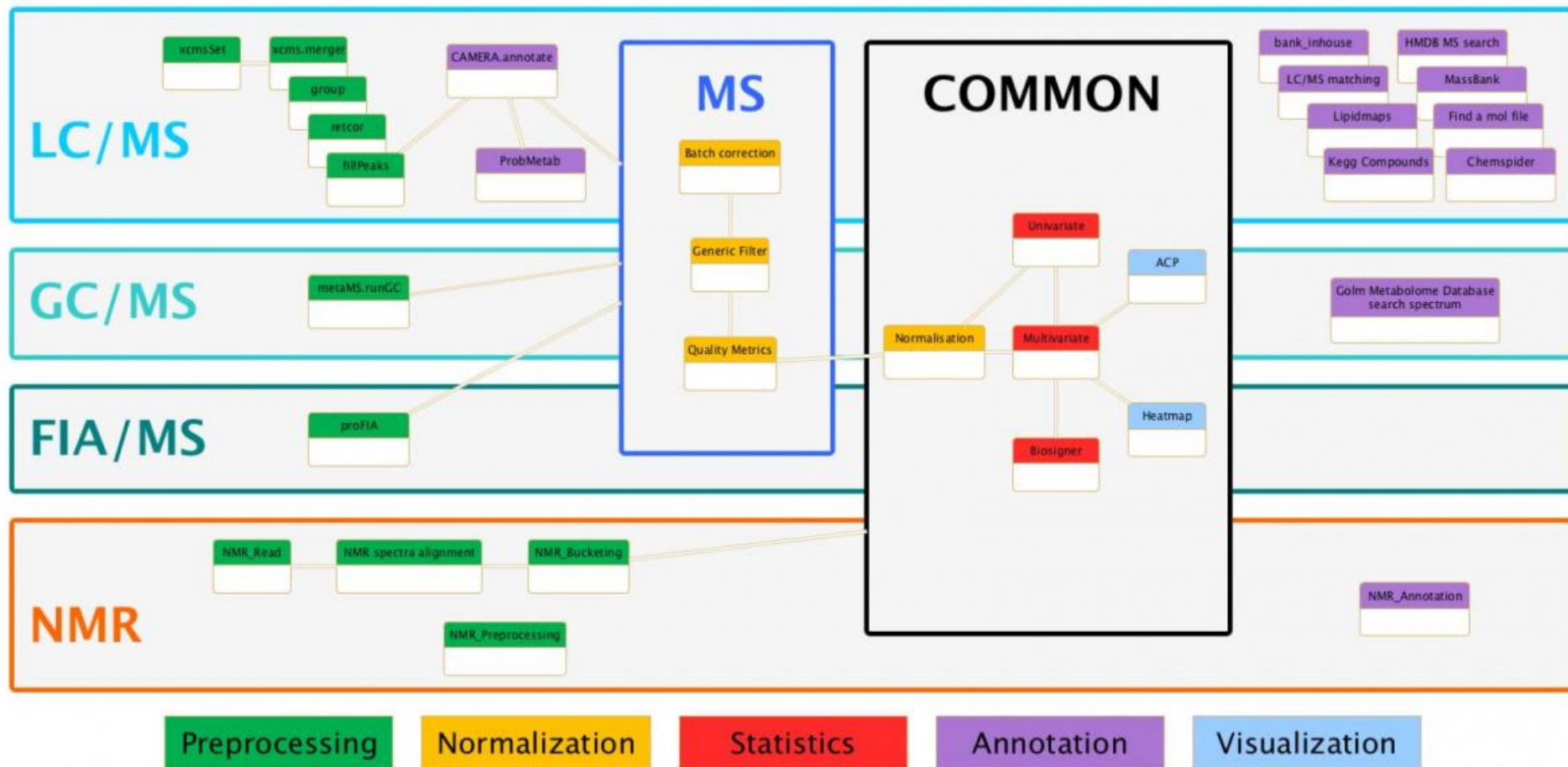
Nutrition



Biotechnology



W4M tools



[Guitton *et al.* \(2017\). Create, run, share, publish, and reference your LC-MS, FIA-MS, GC-MS, and NMR data analysis workflows with the Workflow4Metabolomics 3.0 Galaxy online infrastructure for metabolomics. *Int. J. Biochem. Cell. Biol.* 93:89-101.](#)

W4M00002_sacurine-comprehensive

Galaxy / 4 / Metabolomics

Analyze Data Workflow Shared Data Visualization Help User

Workflow Canvas | W4M00002w_nitbis2

Preprocessing

Input dataset x xcms.xcmsSet x xcms.group x xcms.retcor x xcms.group x xcms.filePeaks x

output Zip file xset RData file xset RData file xset RData file xset RData file xset RData file

sampleMetadata (tabular) xset RData (rdata.xcms.group) xset RData (rdata.xcms.retcor) xset RData (rdata.xcms.group) xset RData (rdata.xcms.filepeaks)

bcsRawPdf (pdf) rplotsPdf (pdf) rplotsPdf (pdf) rplotsPdf (pdf) rplotsPdf (pdf)

bpcsRawPdf (pdf) log (txt) ticsCorPdf (pdf) bpcsCorPdf (pdf)

Statistics

Transformation x Quality Metrics x Univariate x Multivariate x

Data matrix file Data matrix file Data matrix file Data matrix file

sampleMetadata_out (tabular) Sample metadata file Sample metadata file Sample metadata file

variableMetadata_out (tabular) Variable metadata file Variable metadata file Variable metadata file

information (txt) figure (pdf) information (txt) information (txt)

information (txt) information (txt) information (txt) information (txt)

Annotation

Generic_filter x MassBank x Compute x Kegg Compounds x

Data Matrix file File of masses (Variable metadata) as a new column to Kegg Compounds

Sample metadata file variableMetadata (tabular) out_file1 File of masses (Variable metadata)

Variable metadata file massBankResView (html) keggResView (html)

dataMatrix_out (tabular) sampleMetadata_out (tabular) variableMetadata_out (tabular)

variableMetadata_out (tabular) massBankResXis (tabular)

Canvas

CAMERA.annotate x

RData file

variableMetadata (tabular)

datamatrix (tabular)

rdata (rdata.camera.quick, rdata.camera.positive, rdata.camera.negative)

output_zip (zip)

Param.

Multivariate

PCA, PLS and OPLS (Galaxy Tool Version 2.2.4)

Data matrix file

Data input 'dataMatrix_in' (tabular)

variable x sample, decimal: '.', missing: NA, mode: numerical, sep: tabular

Sample metadata file

Data input 'sampleMetadata_in' (tabular)

sample x metadata, decimal: '.', missing: NA, mode: character and numerical, sep: tabular

Variable metadata file

Data input 'variableMetadata_in' (tabular)

variable x metadata, decimal: '.', missing: NA, mode: character and numerical, sep: tabular

Y Response (for PLS(-DA) and OPLS(-DA) only)

class

Notes: 1) PCA: keep the default (none); 2) PLS(-DA) and OPLS(-DA): indicate the name of the column of the sample table to be modeled

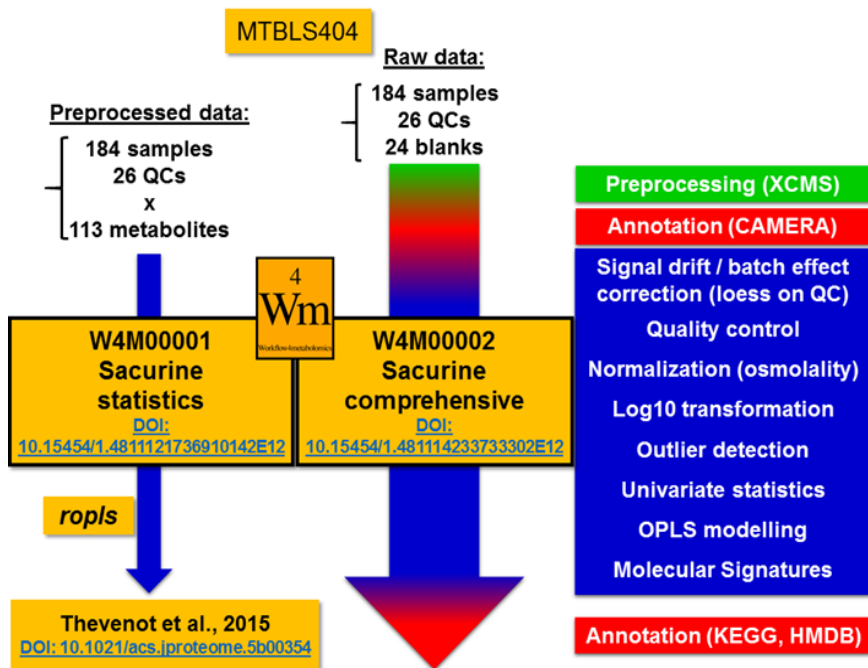
Number of predictive components

NA

Notes: 1) PCA and PLS(-DA): NA can be selected to get a

cea Referencing our analyses

- Demonstrate the value and the reproducibility of your analysis (e.g., to reviewers)
- Receive feedback on your results, get cited, and initiate new collaborations



Referenced W4M histories

WOI	Name & DOI	Technology	Species	Matrice	Factor	Samples
W4M00001	"Sacurine-statistics" 10.15454/1.4811121736910142E12	LC-MS	<i>H. sapiens</i>	Urine	age, BMI, gender	184
W4M00002	"Sacurine-comprehensive" 10.15454/1.481114233733302E12	LC-MS	<i>H. sapiens</i>	Urine		9
W4M00003	"Diaplasma" 10.15454/1.4811165052113186E12	LC-MS	<i>H. sapiens</i>	Plasma		10
W4M00004	"GCMS Algae" 10.15454/1.4811272313071519E12	GC-MS	<i>E. siliculosus</i>	Algae		1
W4M00005	"Ractopamine" 10.15454/1.4811287270056958E12	LC-MS	<i>S. scrofa</i>	Serum		11
W4M00006	"BPA-MMusculusus" 10.15454/1.4821558812795176E12	NMR	<i>M. musculus</i>	Brain		12
W4M00007	"Coffea leaves" 10.15454/1.4985472277740251E12	LCMS	<i>Coffea sp.</i>	Leaves		

[Guitton et al. \(2017\). Create, run, share, publish, and reference your LC-MS, FIA-MS, GC-MS, and NMR data analysis workflows with the Workflow4Metabolomics 3.0 Galaxy online infrastructure for metabolomics. *Int. J. Biochem. Cell. Biol.* 93:89-101.](#)

W4M offer



The W4M Core Team



- Private account
- Computation and storage resources
- Help desk

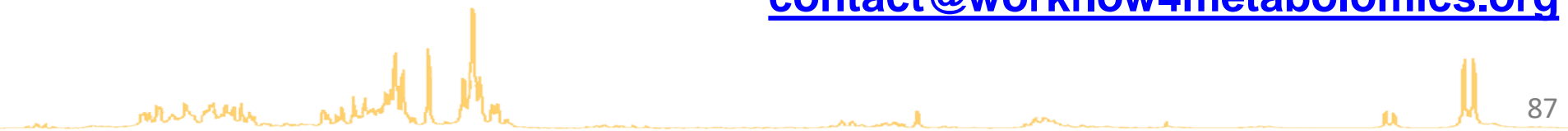
- Sharing and referencing of histories and workflows (DOI)
- Annual courses (tutoring on your own data)

Save the date: 8-12 October 2018, Pasteur Institute (Paris)

- Installation of local instances



contact@workflow4metabolomics.org



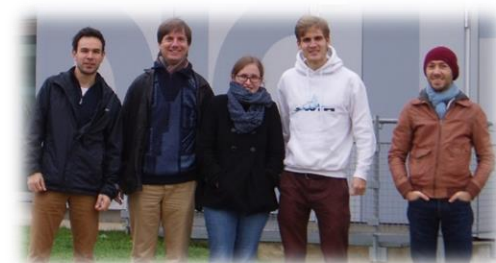
➤ **The Team**



A team work from talented people

➤ The CEA team

- Natacha Lenuzza, Alexis Delabrière, Pierrick Roger-Mele, Bertrand Monfort, Philippe Rinaudo, *et al.*



➤ The MetaboHUB infrastructure



- Fabien Jourdan, Franck Giacomoni, Marie Tremblay-Franco, Jean-François Martin, Mélanie Pétéra, Nils Paulhe, Christophe Junot, Estelle Pujos-Guillot, Dominique Rolin, *et al.*

➤ The IFB infrastucture



- Christophe Caron, Gildas Le Corguillé, David Vallenet, Claudine Médigue, Jacques van Helden, *et al.*

➤ The PhenoMeNal consortium



- Christoph Steinbeck, Steffen Neumann, Namrata Kale, Pablo Moreno, Kenneth Haug, Reza Salek, Philippe Rocca-Serra, Luca Pirredu, *et al.*

Thanks for coming

Questions?