



3^e Journée Intégrative de Protéomique et Métabolomique (8 octobre 2020 en visioconférence)



DE LA RECHERCHE À L'INDUSTRIE

- ▶ Data Sciences for Molecular Phenotyping and Precision Medicine team
 - ▶ CEA, INRAE, Paris Saclay University, MetaboHUB, 91191 Gif-sur-Yvette, France
 - ▶ <https://scidophenia.github.io> SciDophenia
- etienne.thevenot@cea.fr

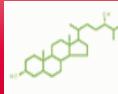
Proteomics and metabolomics data integration: Where do we stand?

What?

Proteomics



Metabolomics



- ▶ large-scale study of proteins
- ▶ post-translational modifications

- ▶ small molecule substrates, intermediates, and products of metabolism
- ▶ peptides, carbohydrates, lipids, nucleosides
- ▶ “functional readout of the physiological state”

Why?

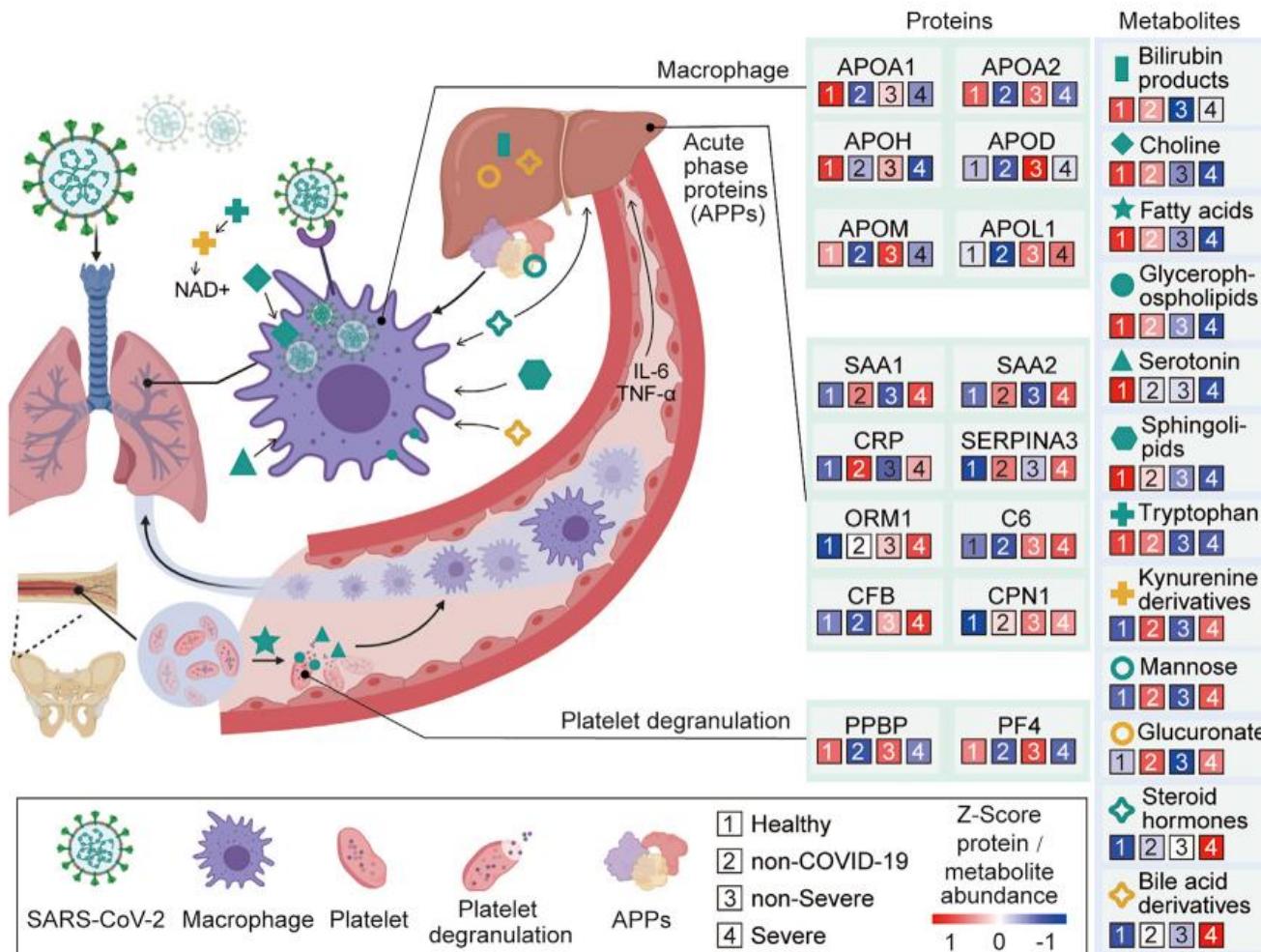
Proteins – Metabolites

► Interactions

- Building blocks of proteins
- Substrates, cofactors, products of enzymatic reactions
- Allosteric regulators (enzymes, receptors, transcription factors)
- Post-translational modifications by covalent link to metabolites

Piazza *et al.* (2018). A map of protein-metabolite interactions reveals principles of chemical communication. *Cell*, **172**:358–372.

Increase the biological understanding

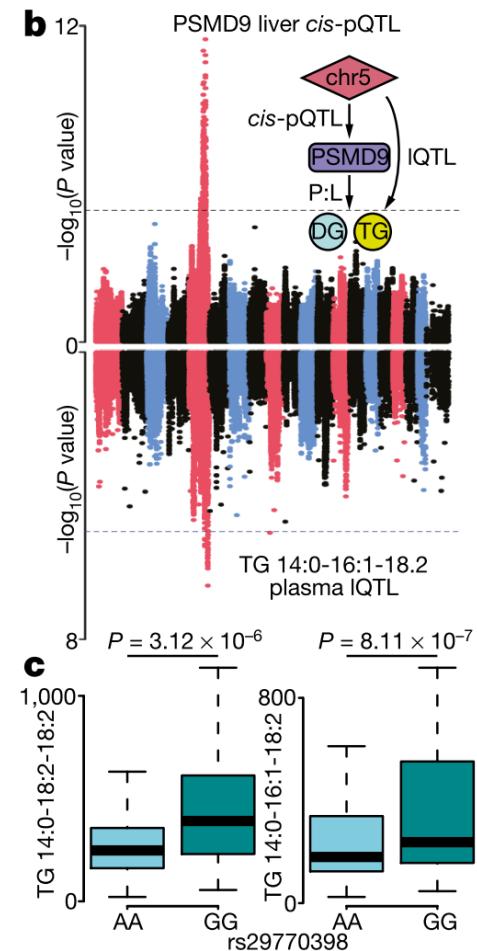
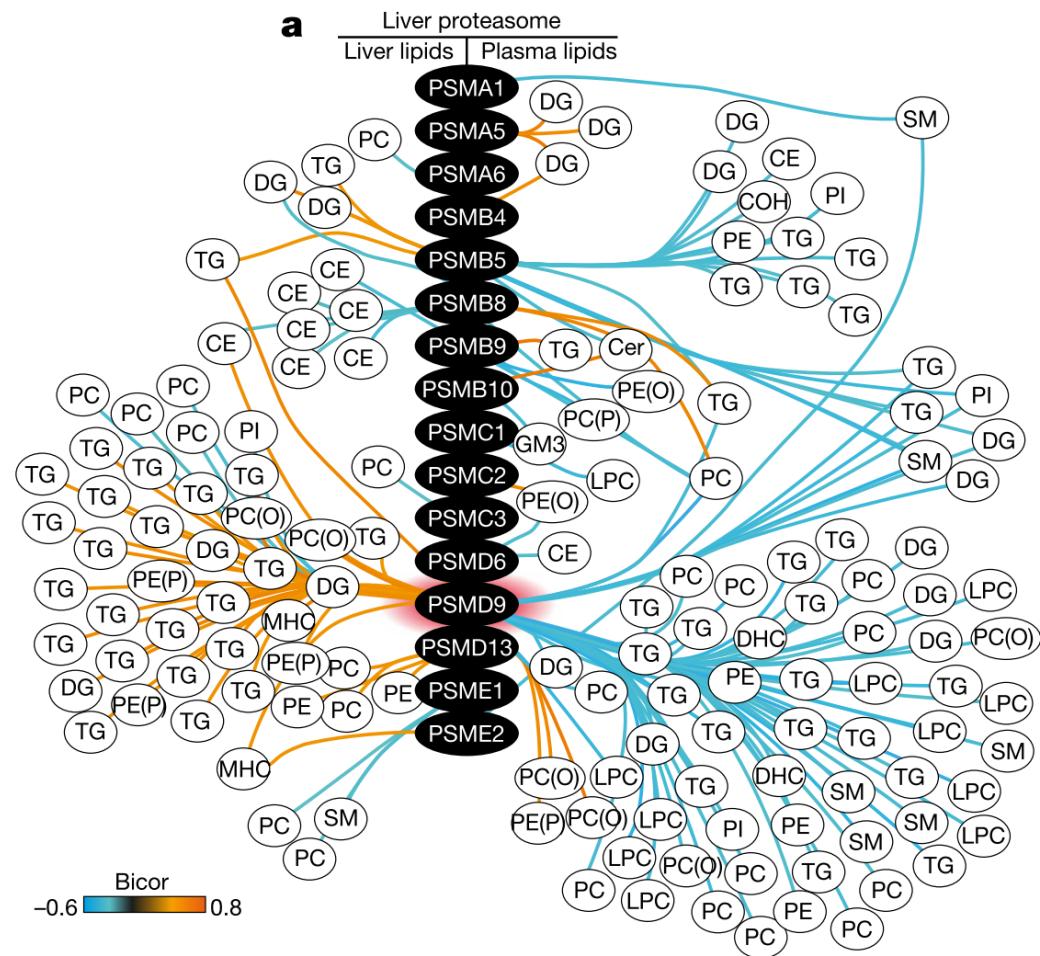
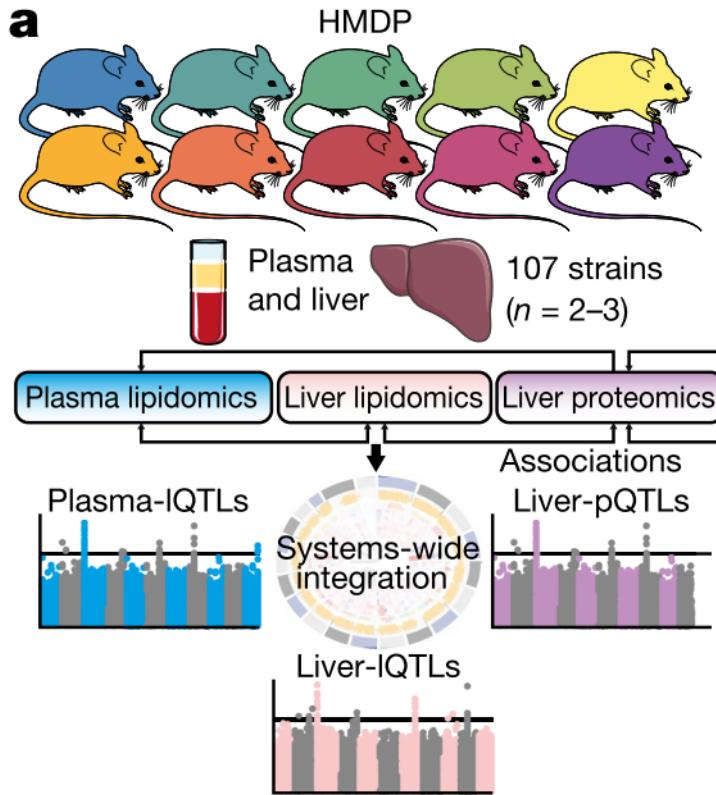


Shen *et al.* (2020). Proteomic and metabolomic characterization of COVID-19 patient sera. *Cell*, 9:59–72.

Figure 5. Key Proteins and Metabolites Characterized in Severe COVID-19 Patients in a Working Model

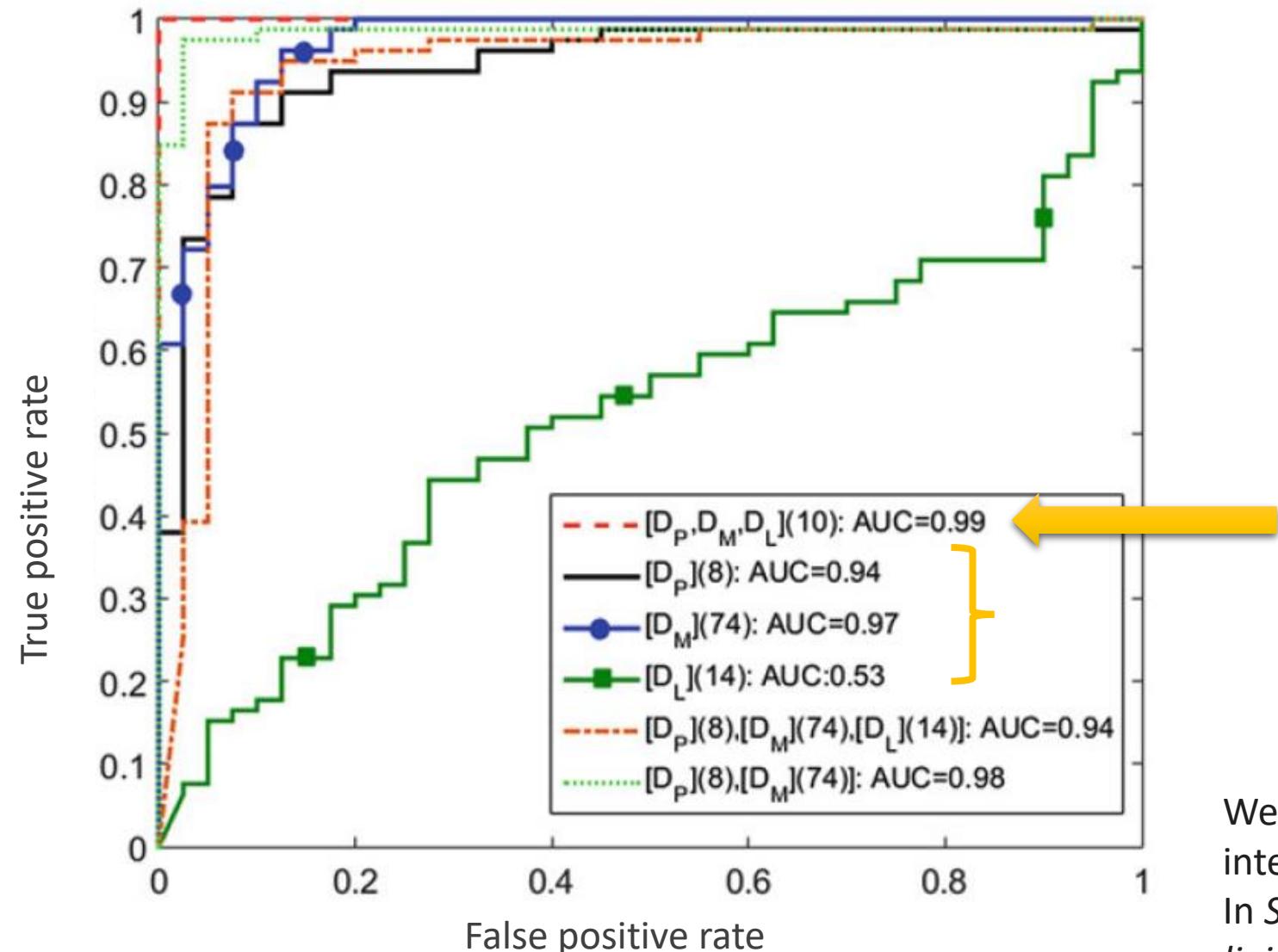
SARS-CoV-2 may target alveolar macrophages via ACE2 receptor, leading to an increase of secretion of cytokines including IL-6 and TNF- α , which subsequently induce the elevation of various APPs such as SAP, CRP, SAA1, SAA2, and C6, which are significantly upregulated in the severe group. Proteins involved in macrophage, lipid metabolism, and platelet degranulation were indicated with their corresponding expression levels in four patient groups.

Elucidate gene function



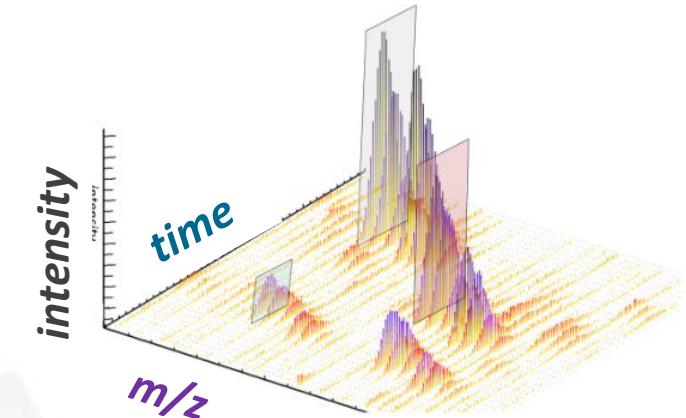
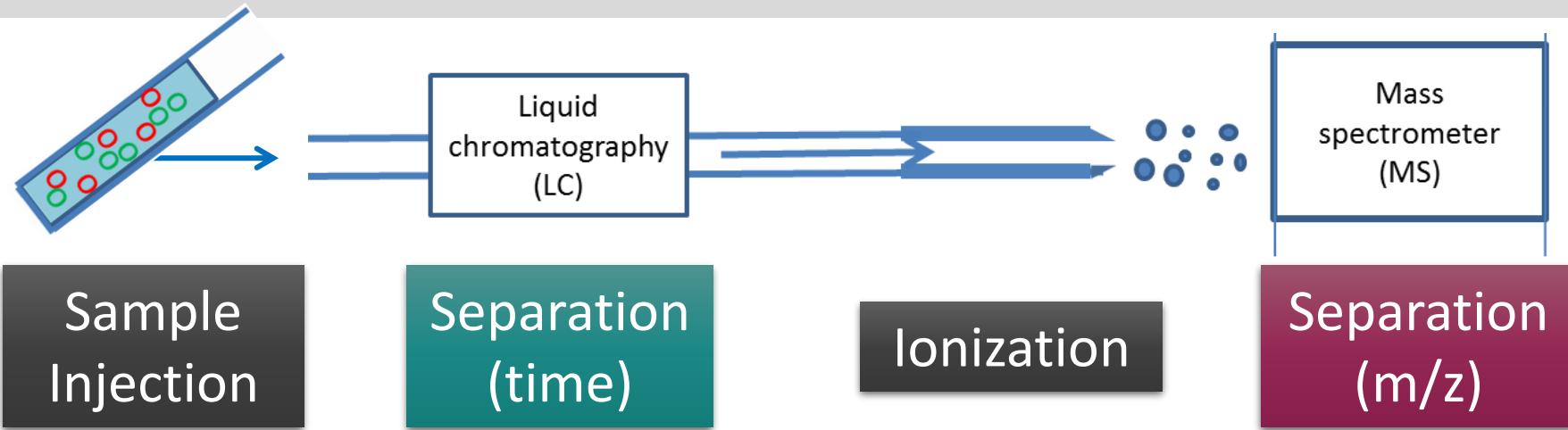
Parker *et al.* (2019). An integrative systems genetic analysis of mammalian lipid metabolism. *Nature*, **567**:187–193.

Increase the predictive performance



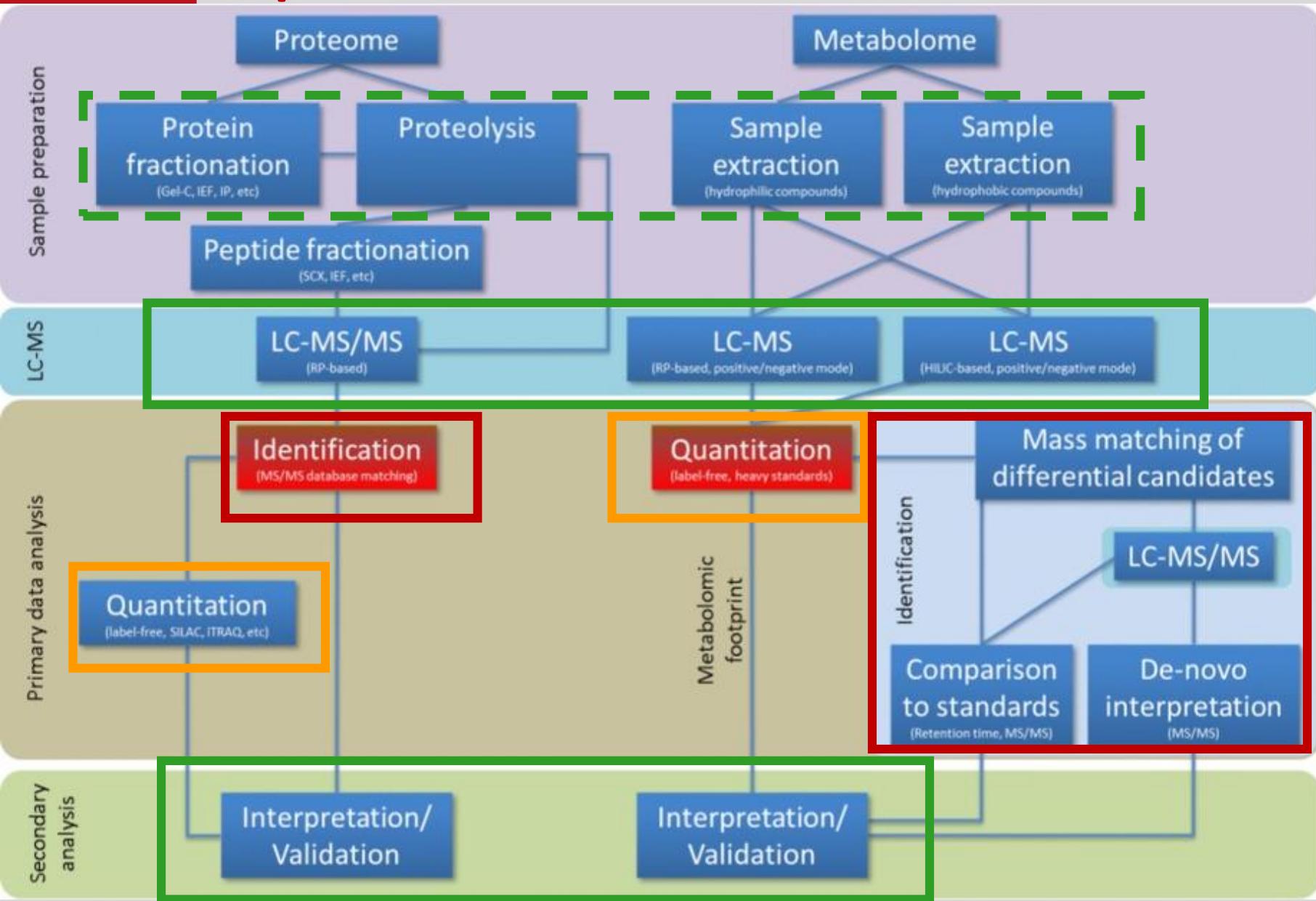
Webb-Robertson *et al.* (2016). Bayesian posterior integration for classification of mass spectrometry data. In *Statistical analysis of proteomics, metabolomics, and lipidomics data using mass spectrometry* (pp. 203–211).

Common technology: LC-HRMS



Similarities in data quantification and statistical analysis

Specificities in data annotation

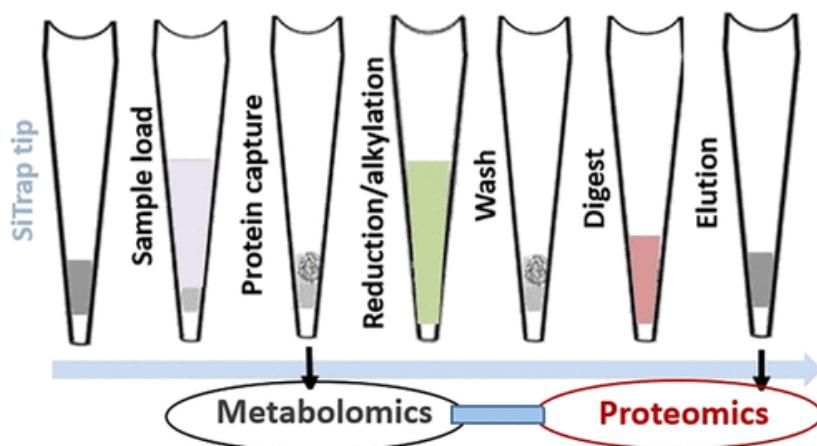


Fischer *et al.* (2013). Two birds with one stone: doing metabolomics with your proteomics kit.
Proteomics, **13**:3371-3386.

Common sample preparation studies

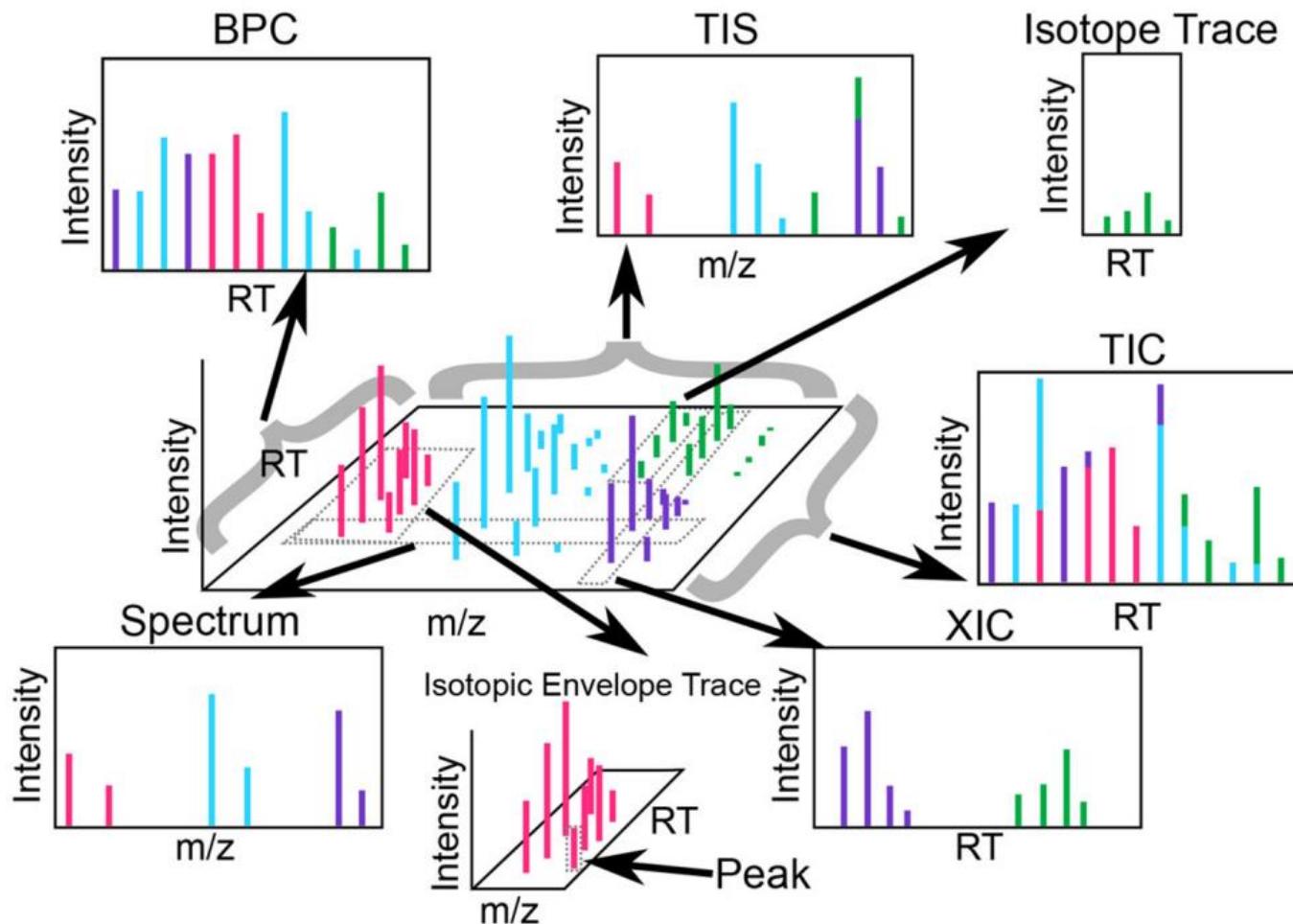
Fischer *et al.* (2013). Two birds with one stone: doing metabolomics with your proteomics kit. *PROTEOMICS*, **13**:3371-3386.

Blum *et al.* (2018). Single-platform ‘multi-omic’ profiling: unified mass spectrometry and computational workflows for integrative proteomics–metabolomics analysis. *Molecular Omics*, **14**:307–319.



Zougman *et al.* (2019). Detergent-free simultaneous sample preparation method for proteomics and metabolomics. *Journal of Proteome Research*, **19**:2838–2844.

Common nomenclature

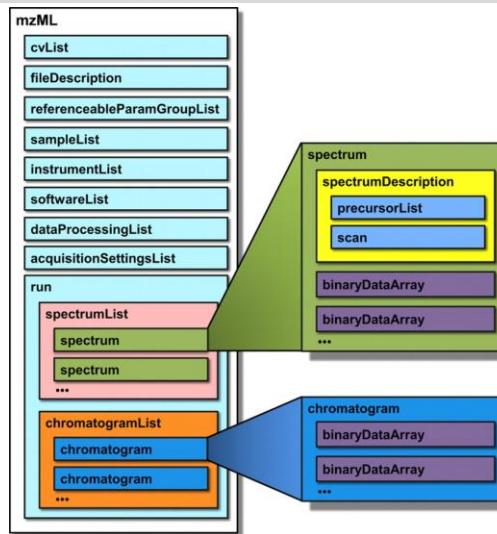


Smith *et al.* (2014). Proteomics, lipidomics, metabolomics: a mass spectrometry tutorial from a computer scientist's point of view. *BMC Bioinformatics*. **15**.

Common formats

► Storage

- raw data: **mzML**
- processed data: **mzTab**
 - quantification
 - identification

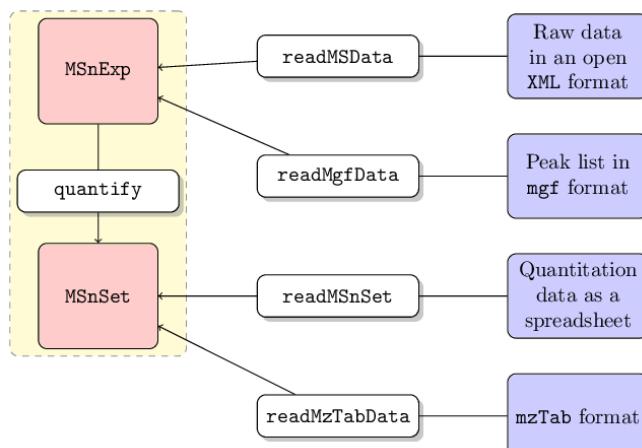


Martens *et al.* (2010). mzML - a community standard for mass spectrometry data. *Molecular & Cellular Proteomics*. **10**.



► Computation

- R object: **MSnbase**



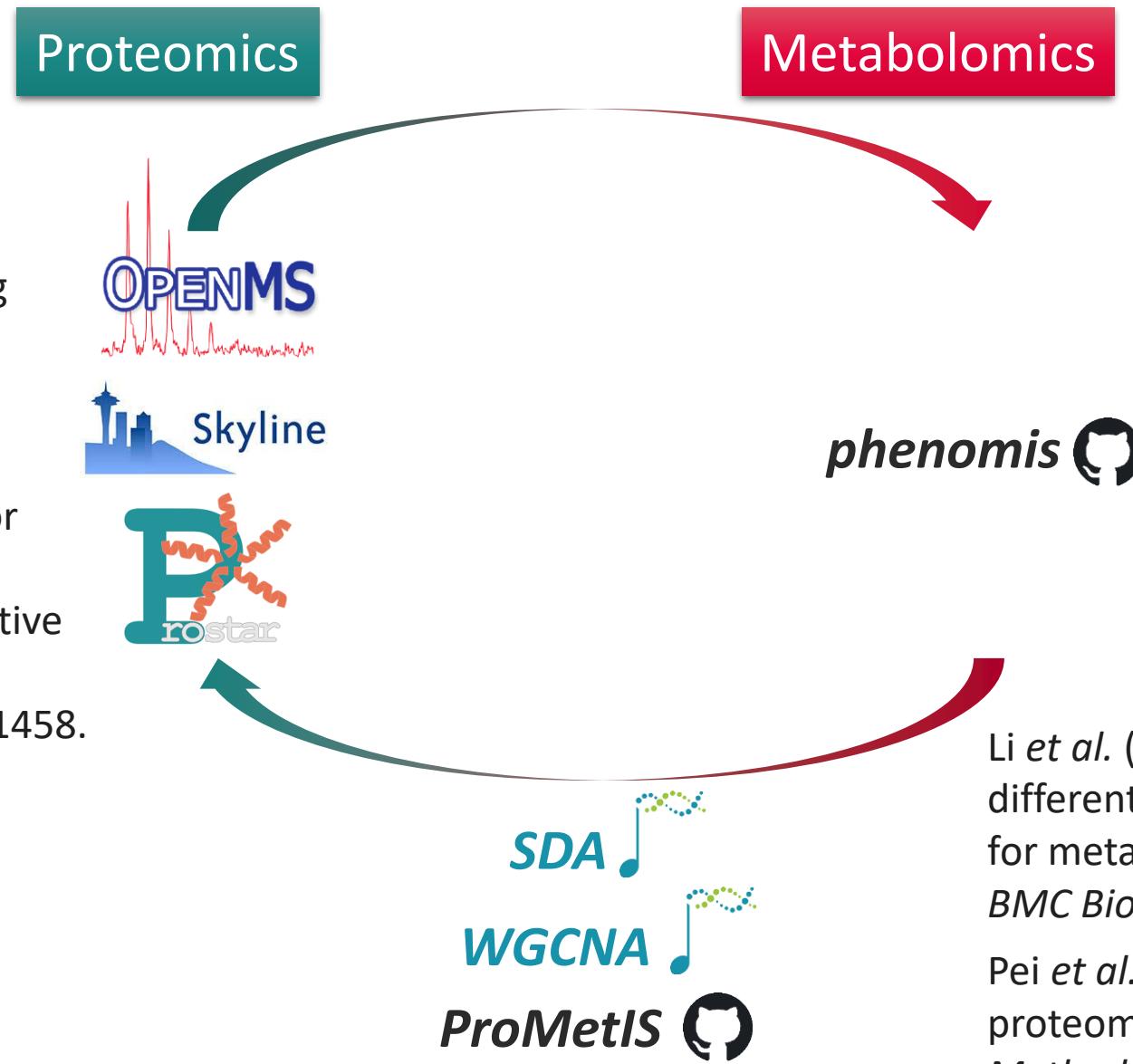
Griss *et al.* (2014). The mzTab data exchange format: Communicating mass-spectrometry-based proteomics and metabolomics experimental results to a wider audience. *Molecular & Cellular Proteomics*. **13**:2765-2775.

Gatto *et al.* (2020). MSnbase, efficient and elegant R-based processing and visualization of raw mass spectrometry data. *Journal of Proteome Research*.

Toward common software

Rurik *et al.* (2020). Metabolomics data processing using OpenMS. *Methods in Molecular Biology*, **2104**.

Adams *et al.* (2020). Skyline for small molecules: A unifying software package for quantitative metabolomics. *Journal of Proteome Research*, **19**:1447-1458.



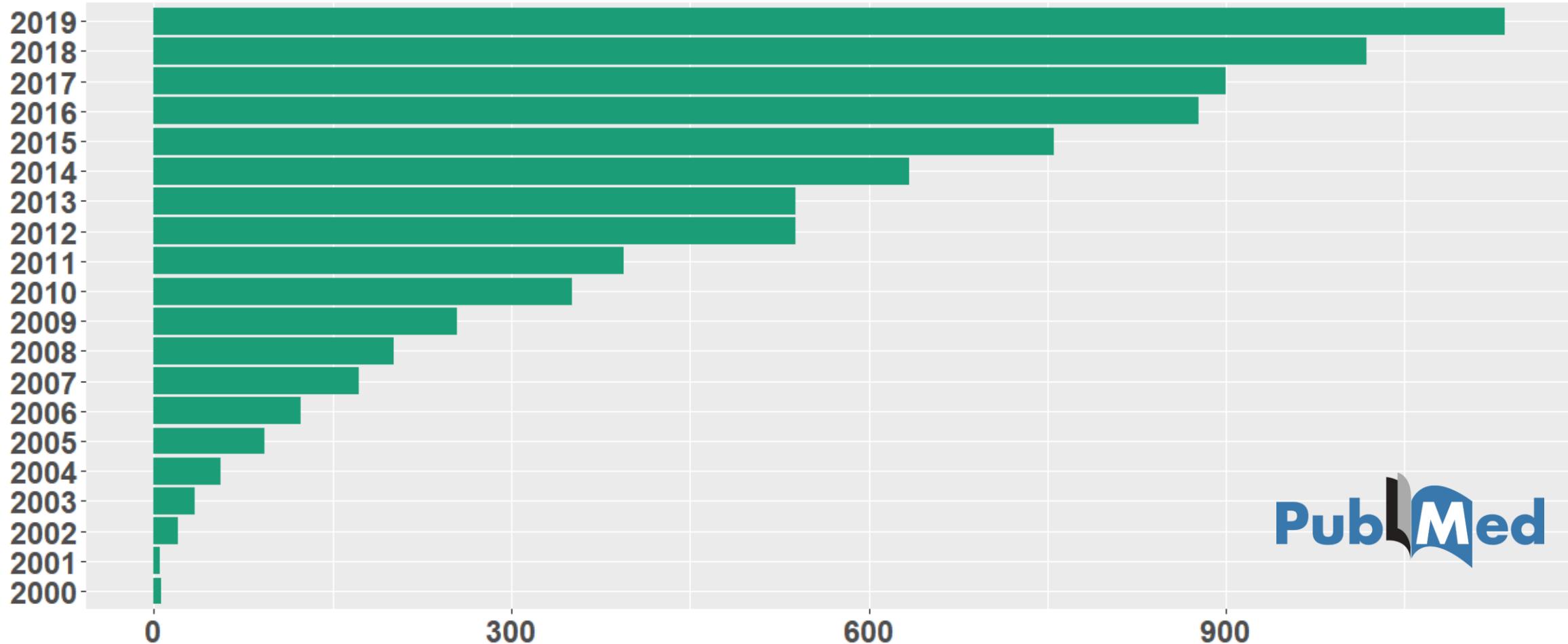
Li *et al.* (2019). SDA: a semi-parametric differential abundance analysis method for metabolomics and proteomics data. *BMC Bioinformatics*. **20**.

Pei *et al.* (2017). WGCNA application to proteomic and metabolomic data analysis. *Methods in Enzymology*. **135**-158.

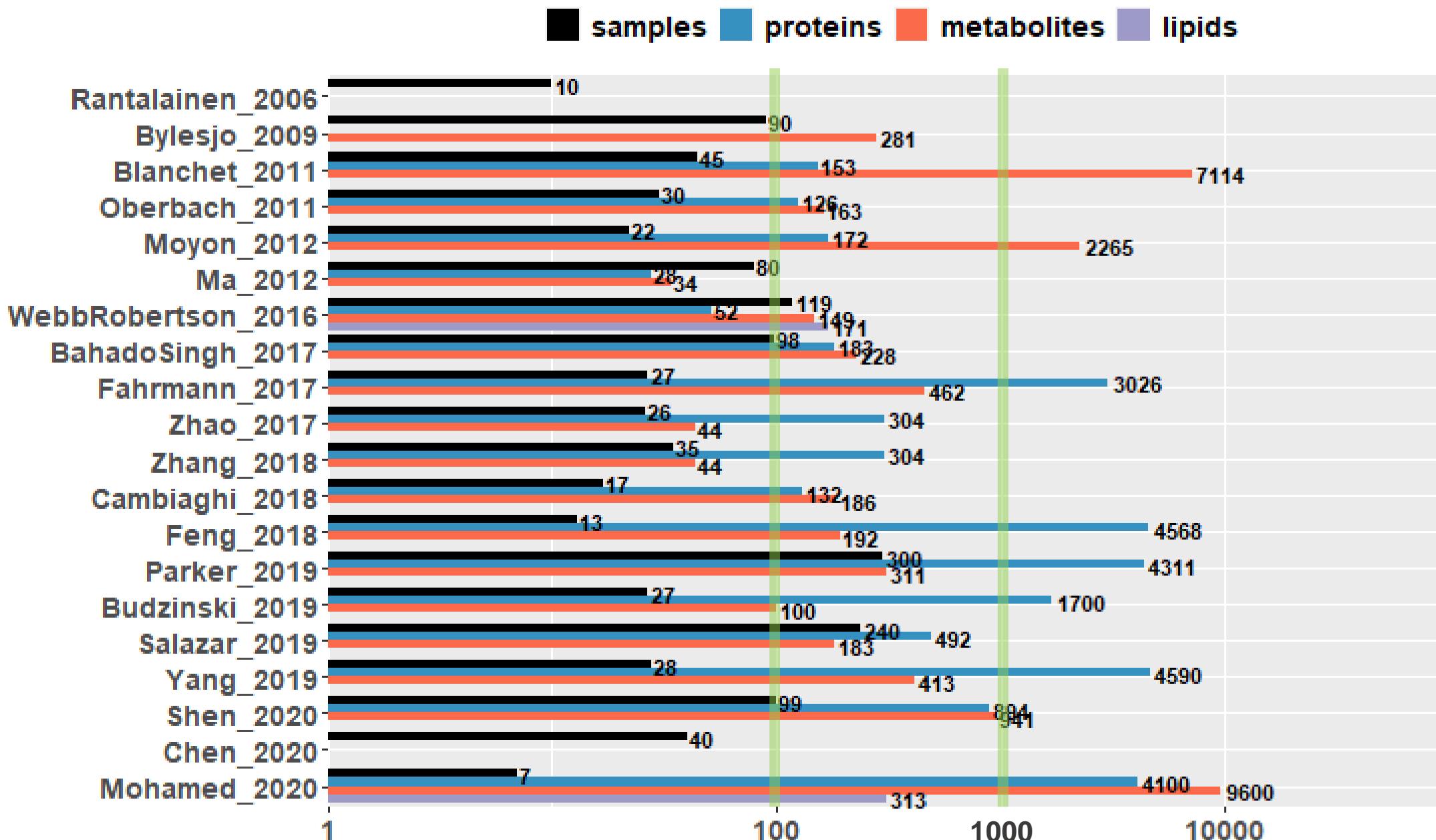
Who?

Commissariat à l'énergie atomique et aux énergies alternatives - www.cea.fr

Joint proteomics & metabolomics publications



Bibliography



► criteria:

- untargeted proteomics and metabolomics
- include a statistical data integration methodology
- open-source implementation

► selection:

- Moyon *et al.* (2012). Statistical strategies for relating metabolomics and proteomics data: a real case study in nutrition research area. *Metabolomics*, **8**, 1090–1101.
- Webb-Robertson *et al.* (2016). Bayesian posterior integration for classification of mass spectrometry data. In *Statistical analysis of proteomics, metabolomics, and lipidomics data using mass spectrometry* (pp. 203–211).
- Cambiaghi *et al.* (2018). An innovative approach for the integration of proteomics and metabolomics data in severe septic shock patients stratified for mortality. *Scientific Reports*, **8**.
- Parker *et al.* (2019). An integrative systems genetic analysis of mammalian lipid metabolism. *Nature*, **567**, 187–193.
- Shen *et al.* (2020). Proteomic and metabolomic characterization of COVID-19 patient sera. *Cell*, **9**, 59–72.

Bibliographic focus: overview

Reference	Disease	Species	Sample	Samples	Clinics	Proteins	Metabolites	Lipids	Data integration
Moyon et al., 2012	maternal diet impact on offspring	rat	plasma, hypothalamus	22		172	2265		PCA, RCCA, MB-PLS
WebbRobertson et al., 2016	diabetes	human	serum	119		52	149	171	concatenation + (RFE +) NB; LDA and RF fusion
Cambiaghi et al., 2018	severe septic shock	human	plasma	17	17	132	186		feature selection (mRMR) + concatenation + classification (elastic net, LDA, PLS-DA)
Parker et al., 2019	hepatosteatosis, metabolic syndrome	mouse	liver, plasma	300		4311	311		clustering based representative feature selection, lasso
Shen et al., 2020	COVID-19	human	serum	99	12	894	941		random forest, mFuzz, IPA

How?

Data integration: questions

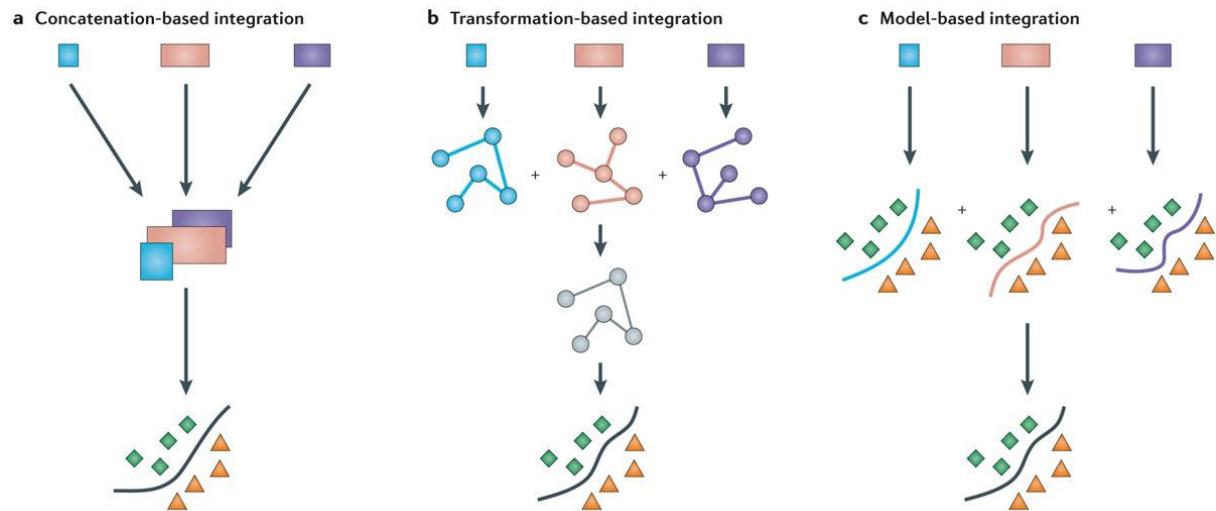
- ▶ Which blocks are the most important for the stratification/prediction?
- ▶ Which features?
- ▶ What is the specific/shared information from each block?
- ▶ How are the features from different blocks correlated?
- ▶ Which biological pathways/networks are significantly involved?

- ▶ **Normalization of each block**
- ▶ **Confounding effects (for each block)**
- ▶ **Overfitting (limited number of samples)**
 - ⇒ validation (statistical, biological)
- ▶ **Feature selection**
- ▶ **Limited annotation of metabolites**
- ▶ **Redundancy/specificity/ambiguity of chemical/biological identifiers in the databases**
- ▶ **Partial coverage of the proteome and metabolome**

Data integration: approaches

► Biostatistics

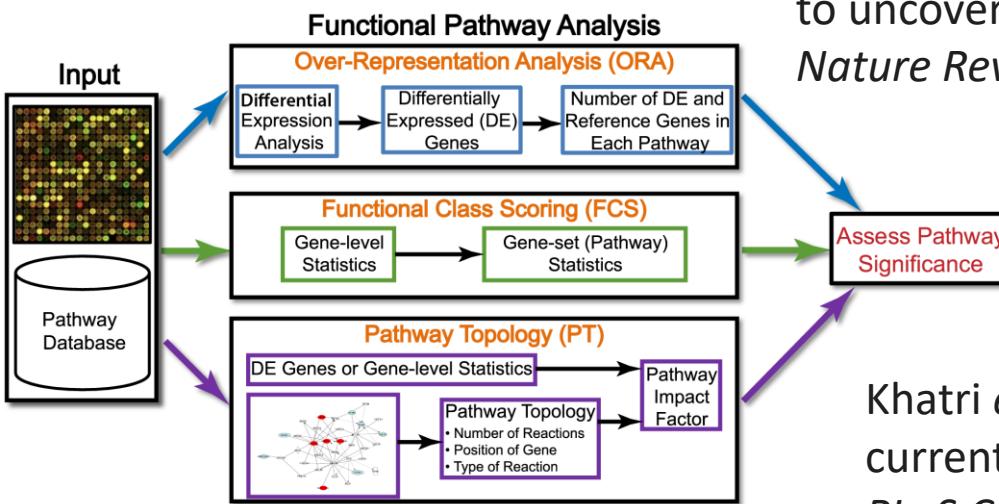
- Fusion
 - low (concatenation of blocks)
 - middle (feature selection/latent variables from each block + model on top)
 - high (one model for each block + vote)
- Correlation networks



Ritchie *et al.* (2015). Methods of integrating data to uncover genotype–phenotype interactions. *Nature Reviews Genetics*, **16**:85–97.

► Bioinformatics

- Mapping
- Enrichment
 - Molecule set
 - Topology-based



Khatri *et al.* (2012). Ten years of pathway analysis: current approaches and outstanding challenges. *PLoS Computational Biology*, **8**: e1002375.

Data integration: (some) published approaches

► Biostatistics

- Fusion
 - low

Unsupervised

- middle

O2PLS (Bylesjö et al., 2009)

- High

Multiblocks PLS (Moyon *et al.*, 2012)

PLS-DA + PCA (Blanchet *et al.*, 2016)

Elastic Net + Elastic Net (Ghaemi *et al.*, 2018)

RGCCA (Budzinski et al., 2019)

LDA, RF (Webb-Robertson *et al.*, 2016)

► Bioinformatics

- Mapping
- Enrichment
 - Molecule set
 - Topology-based

IMPaLA (Bahado-Singh *et al.*, 2017)

ProMetIS: deep phenotyping of mouse models by proteomics and metabolomics

► Objective: high-throughput integration of proteomics and metabolomics data

- innovative methods
- high-quality datasets
- software tools
- workflows

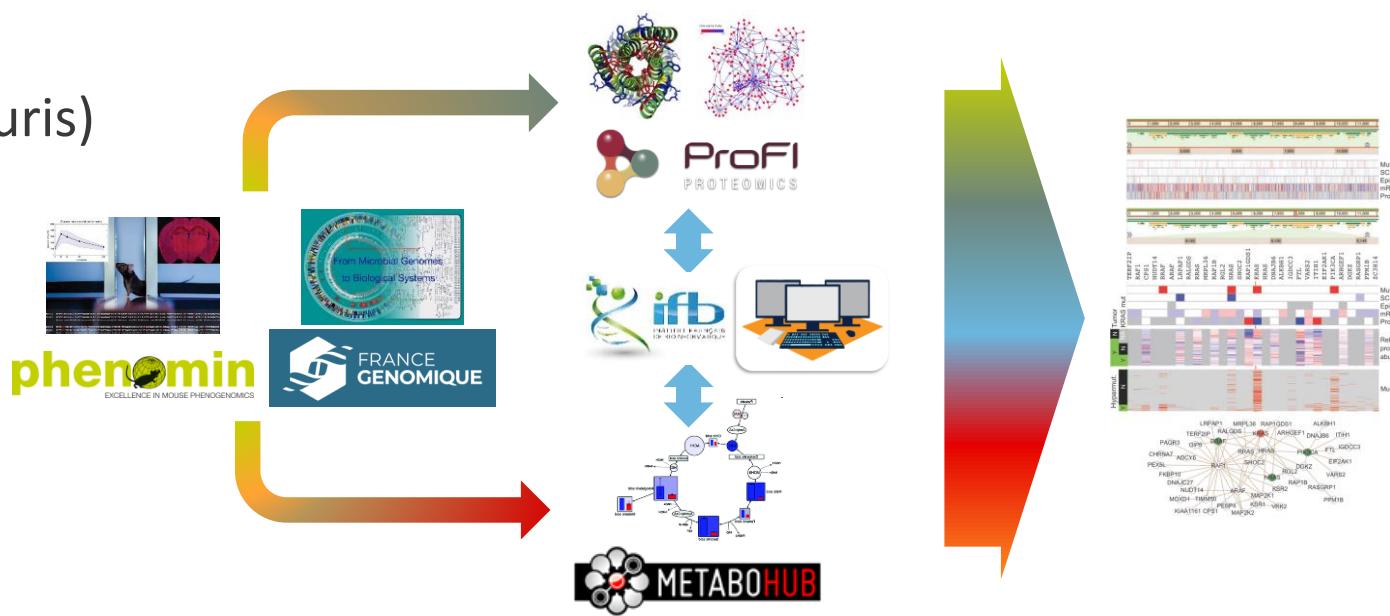
► Case study: molecular phenotyping of mouse models from the IMPC consortium

► Partner infrastructures

- France Génomique
- PHENOMIN (Institut Clinique de la Souris)
- ProFI proteomics
- MetaboHUB
- Institut Français de Bioinformatique

► Post-doctorate

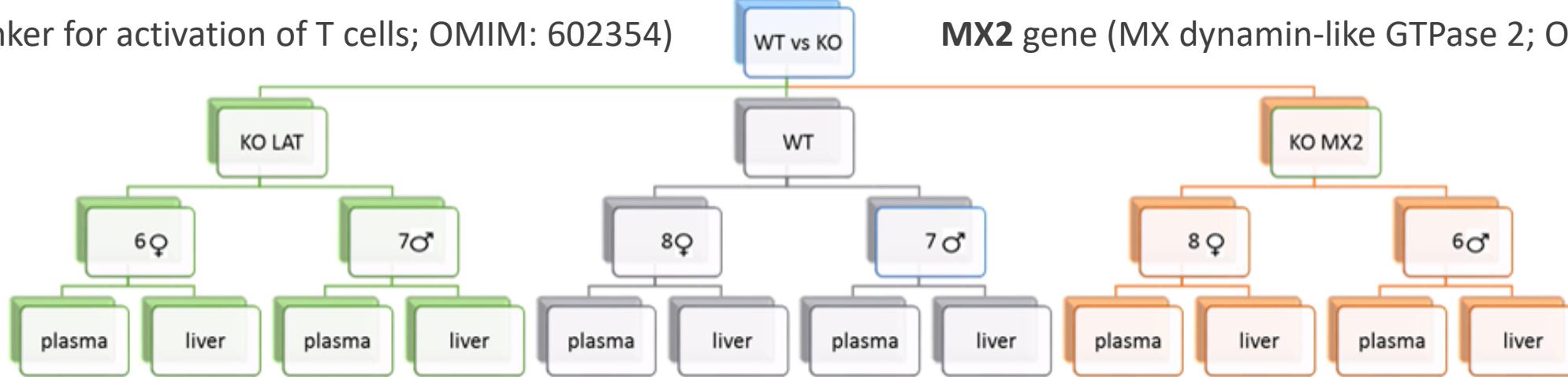
- Alyssa Imbert





phenomin
EXCELLENCE IN MOUSE PHENOMICS

LAT (linker for activation of T cells; OMIM: 602354)



MX2 gene (MX dynamin-like GTPase 2; OMIM: 147890)

► LAT involved in:

- T-cell receptor (TCR) signaling
- Neurodevelopmental diseases

Roncagalli et al. (2010). LAT signaling pathology: an "autoimmune" condition without T cell self-reactivity. *Trends in Immunology*, **31**:253–259.

Loviglio et al. (2017). The immune signaling adaptor LAT contributes to the neuroanatomical phenotype of 16p11.2 BP2-BP3 CNVs. *The American Journal of Human Genetics*, **101**:564–577.



Preclinical



Proteomics

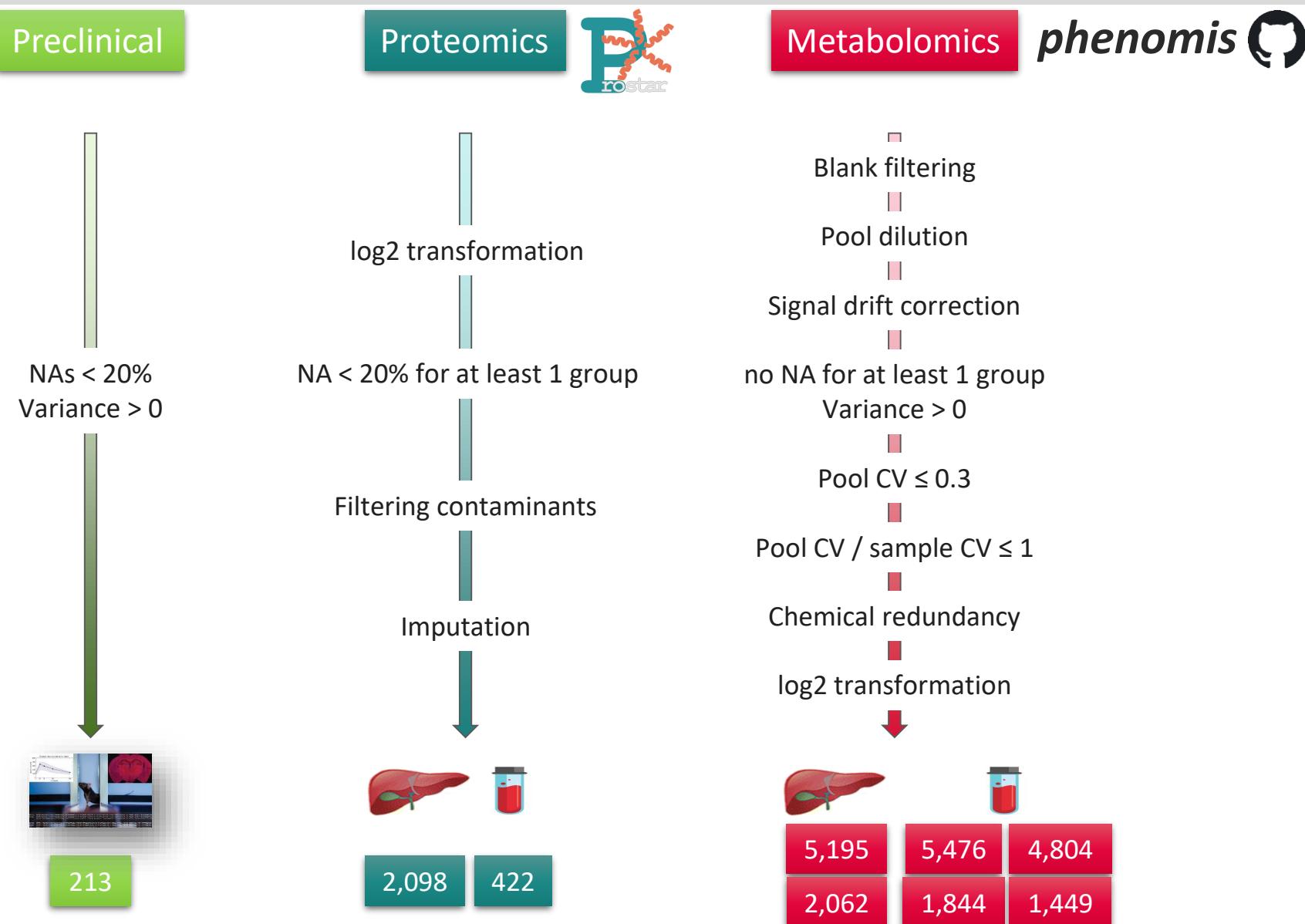


Metabolomics



C18 pos Hypersil	C18 pos Hypersil	C18 pos Acquity
Hilic neg	Hilic neg	C18 neg Acquity

Post-processing



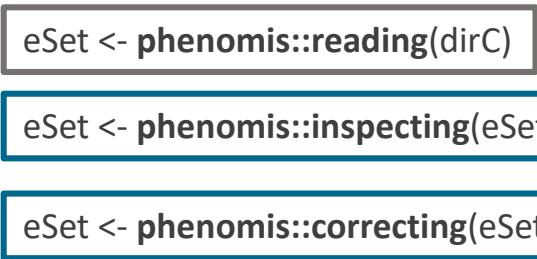
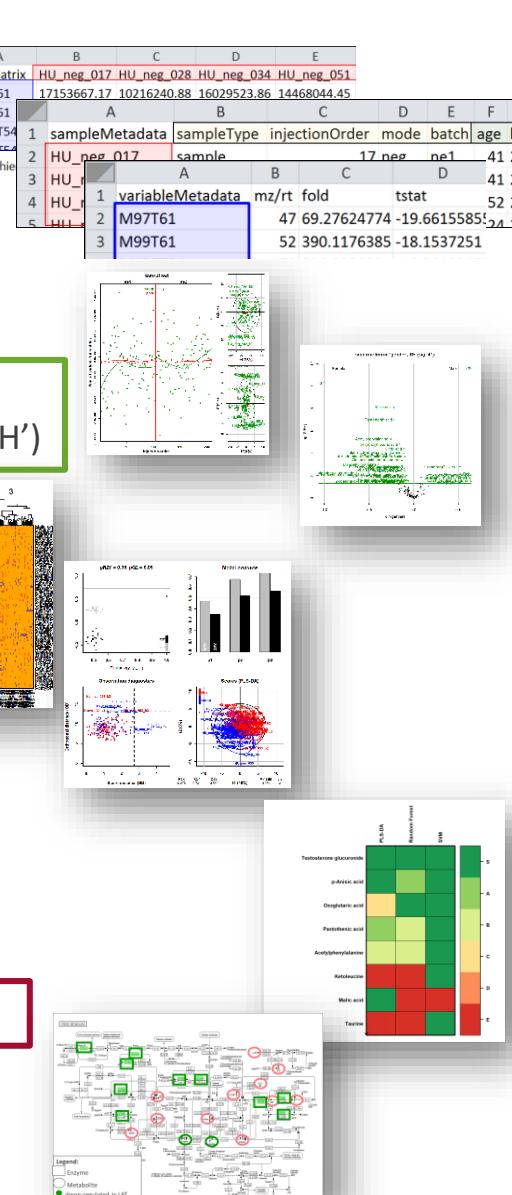
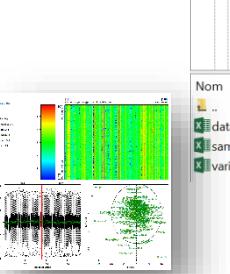
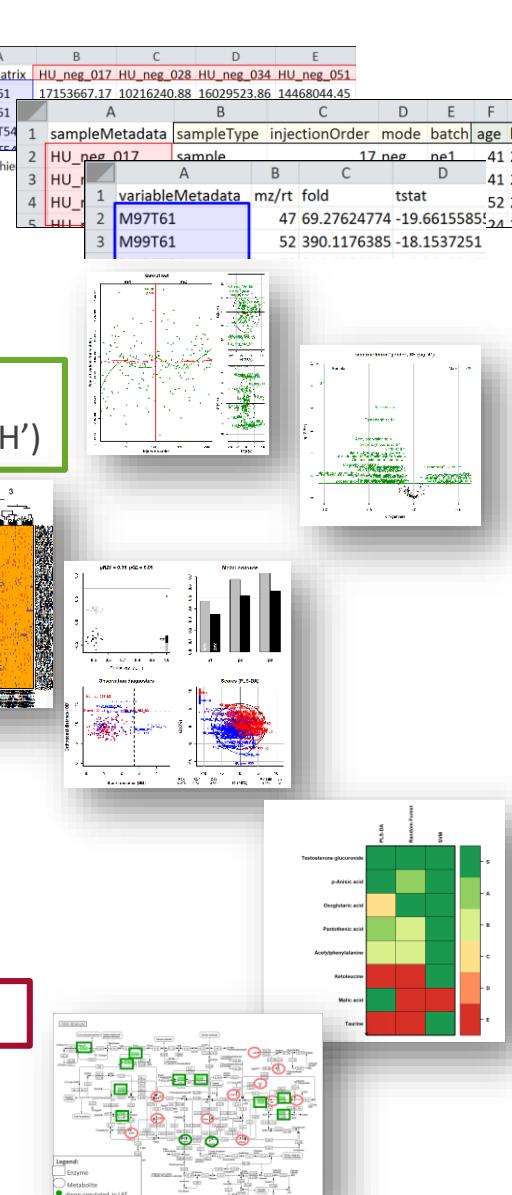
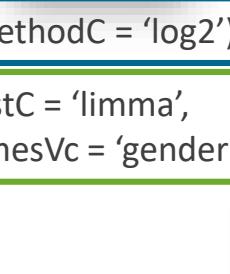
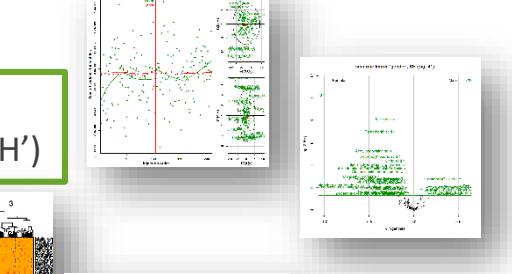
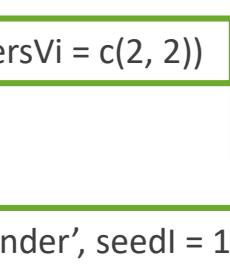
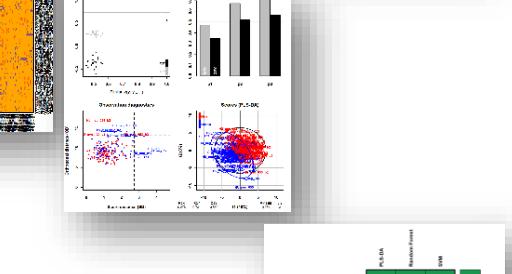
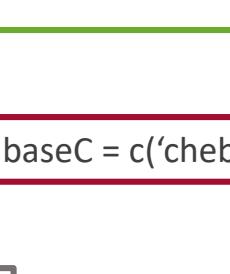
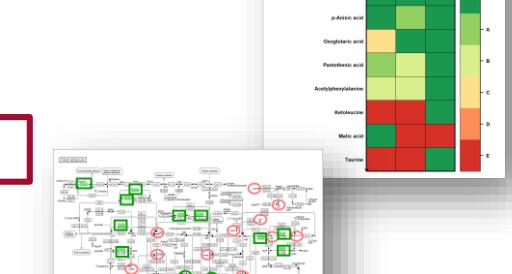
Single omics data analysis workflow

phenomis 
Thévenot et al.

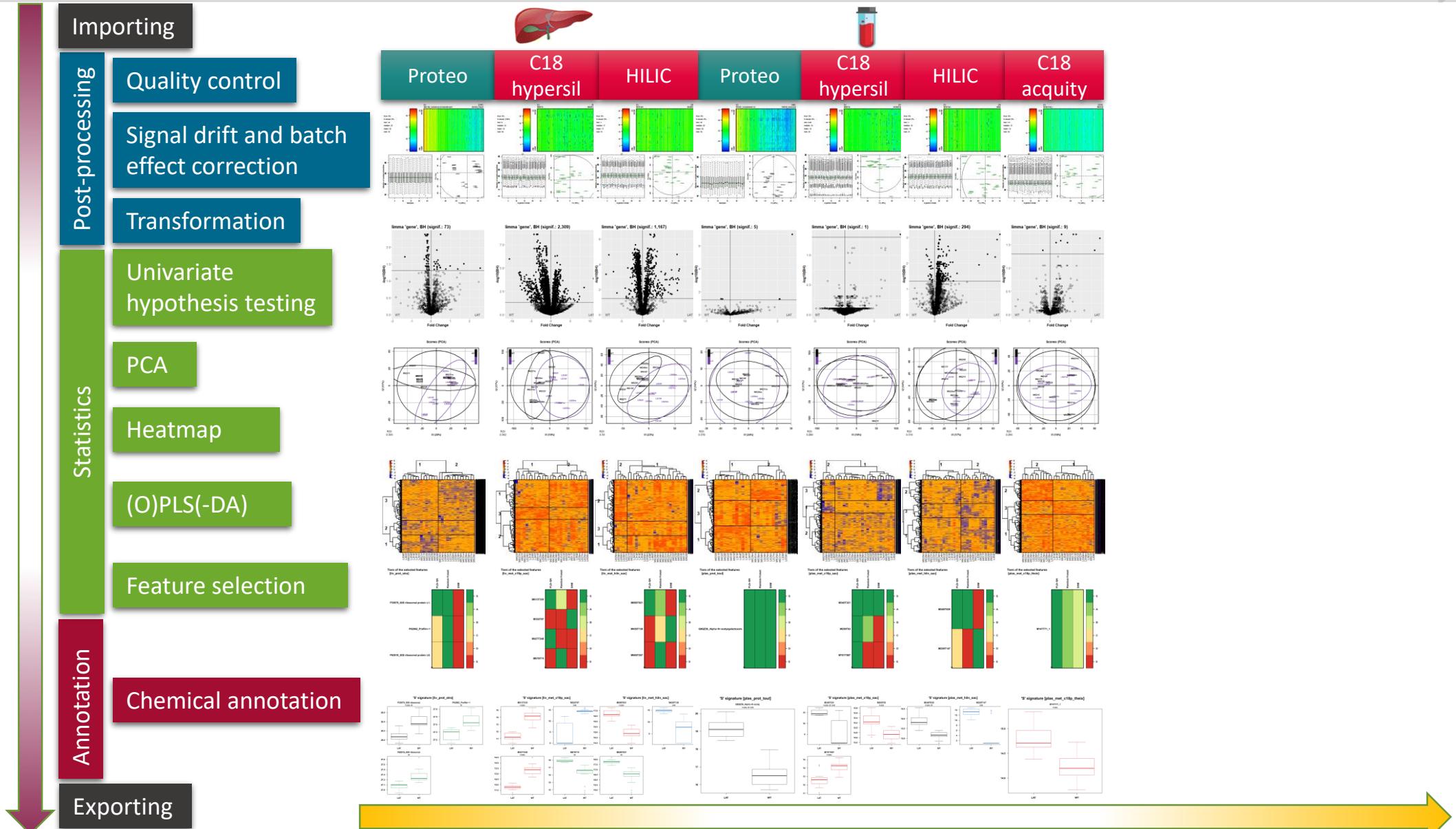
ropls 
Thévenot et al., 2015

biosigner 
Rinaudo et al., 2016

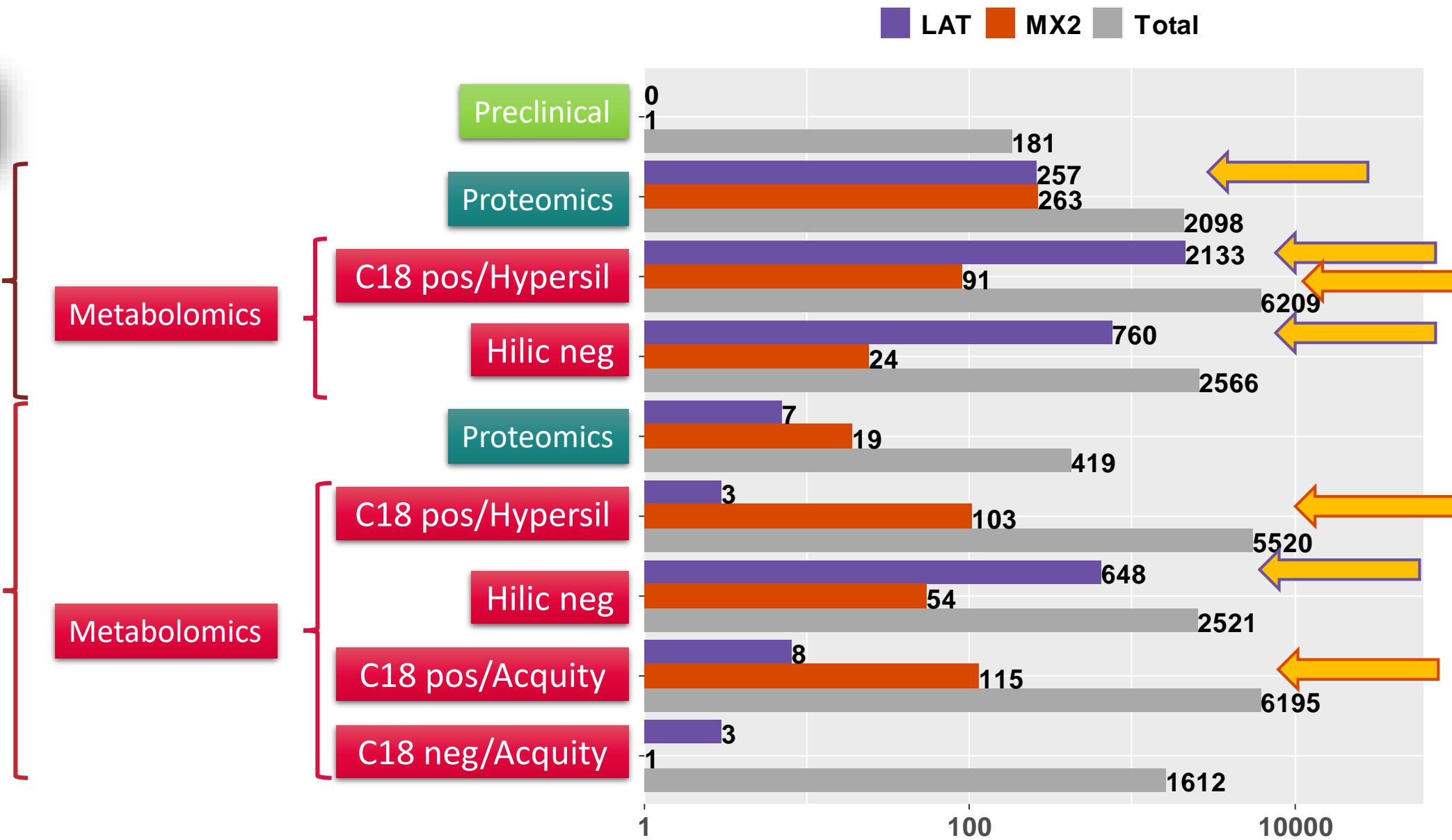
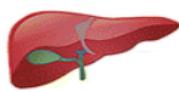
biodb 
Roger et al.

Importing	eSet <- phenomis::reading(dirC)			
Post-processing	Quality control	eSet <- phenomis::inspecting(eSet)		
	Signal drift and batch effect correction	eSet <- phenomis::correcting(eSet)		
	Transformation	eSet <- phenomis::transforming(eSet, methodC = 'log2')		
	Univariate hypothesis testing	eSet <- phenomis::hypotesting(eSet, testC = 'limma', factorNamesVc = 'gender', adjustC = 'BH')		
	PCA	setPca <- ropls::opls(eSet) eSet <- ropls::getEset(setPca)		
	Heatmap	eSet <- phenomis::clustering(eSet, clustersVi = c(2, 2))		
	(O)PLS(-DA)	setPlsda <- ropls::opls(eSet, 'gender') eSet <- ropls::getEset(setPlsda)		
	Feature selection	setBiosign <- biosigner::biosign(eSet, 'gender', seedl = 123) eSet <- biosigner::getEset(setBiosign)		
	Chemical annotation	eSet <- phenomis::annotating(eSet, databaseC = c('chebi', 'local.ms'))		
	Exporting	phenomis::writing(eSet, dirC = getwd())		

Multi-steps and Multi-datasets (platforms, tissues, omics)



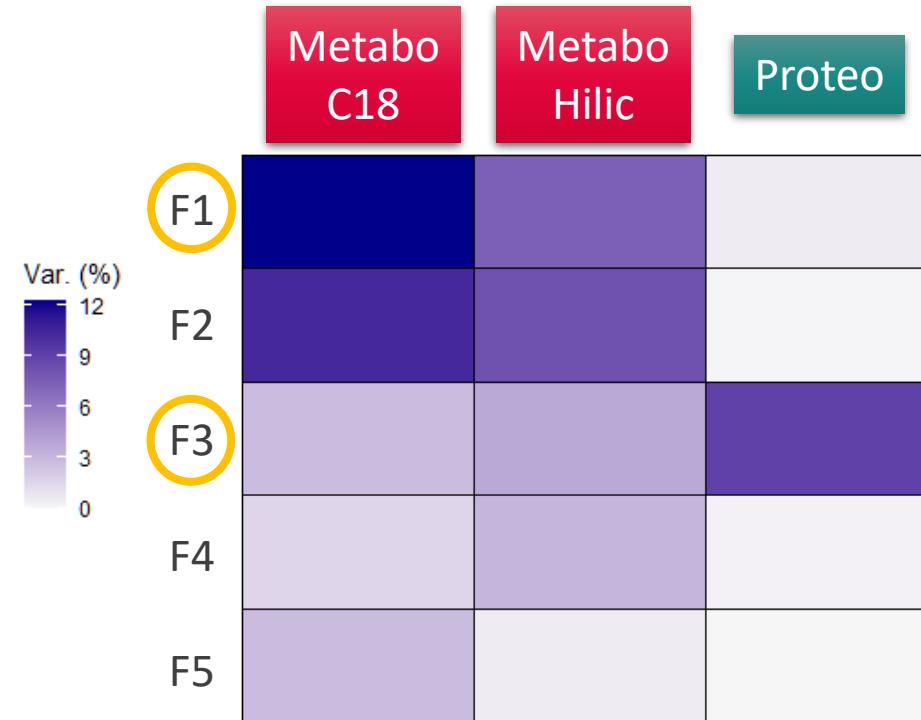
Significant features KO vs WT (limma)



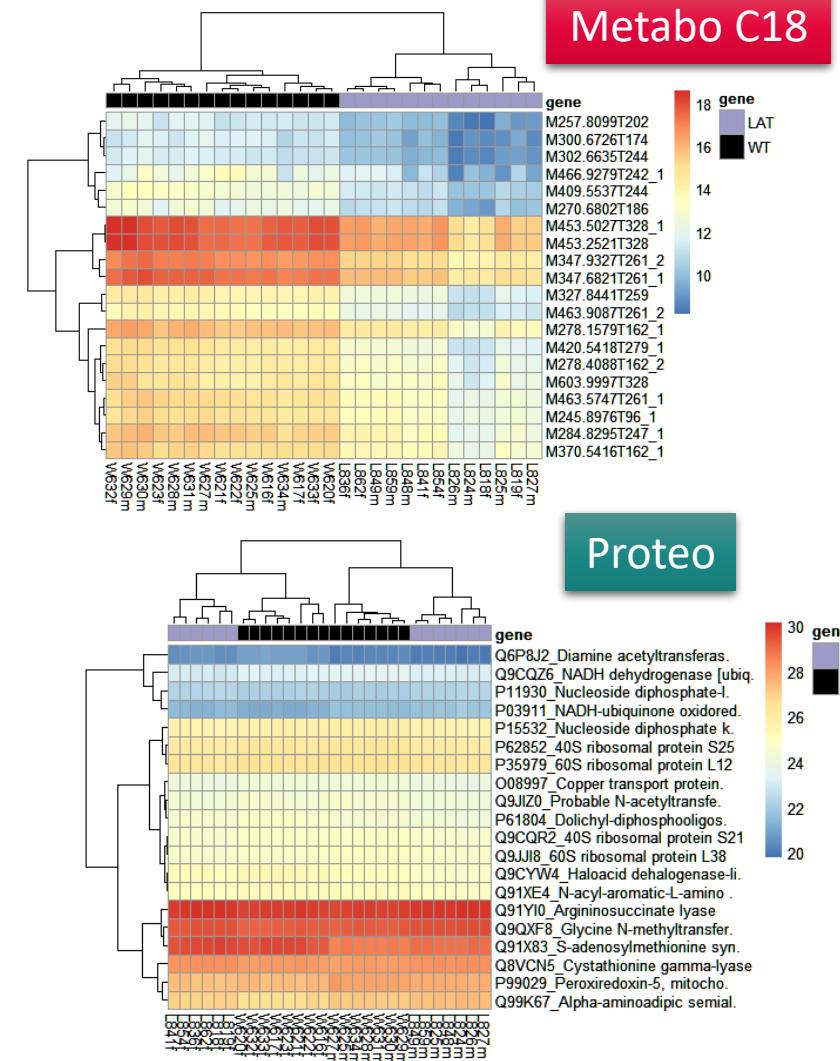
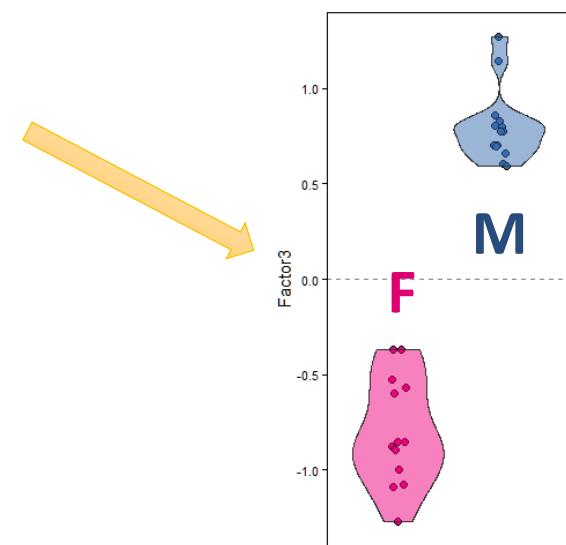
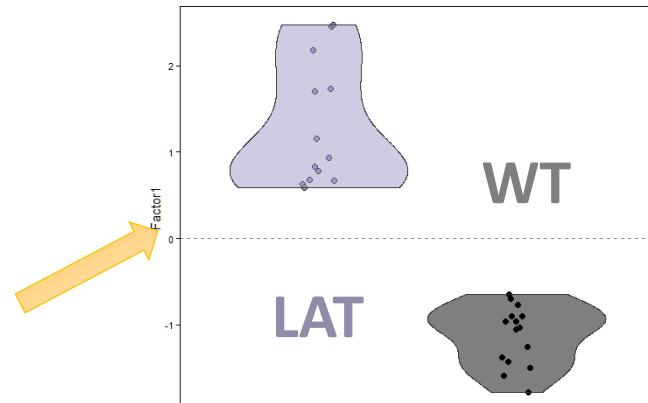
Multi-Omics Factor Analysis



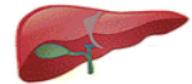
LAT vs WT



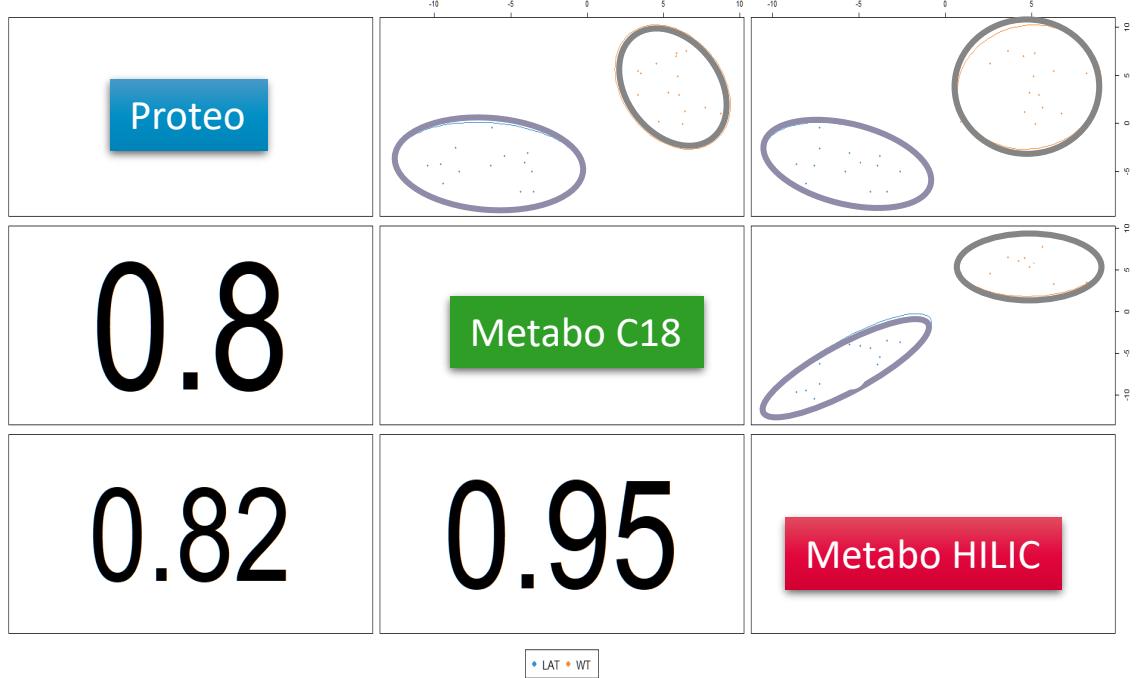
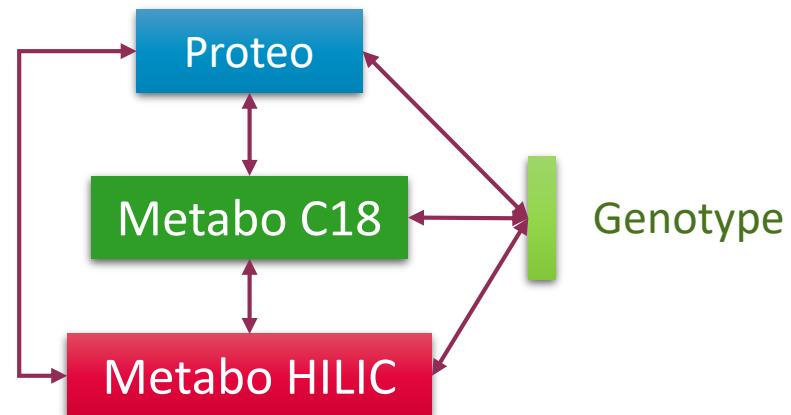
Argelaguet *et al.* (2018). Multi-Omics Factor Analysis—a framework for unsupervised integration of multi-omics data sets. *Molecular Systems Biology*, 14.



Sparse Generalized Canonical Correlation Analysis - *ProMetIS* Discriminant Analysis

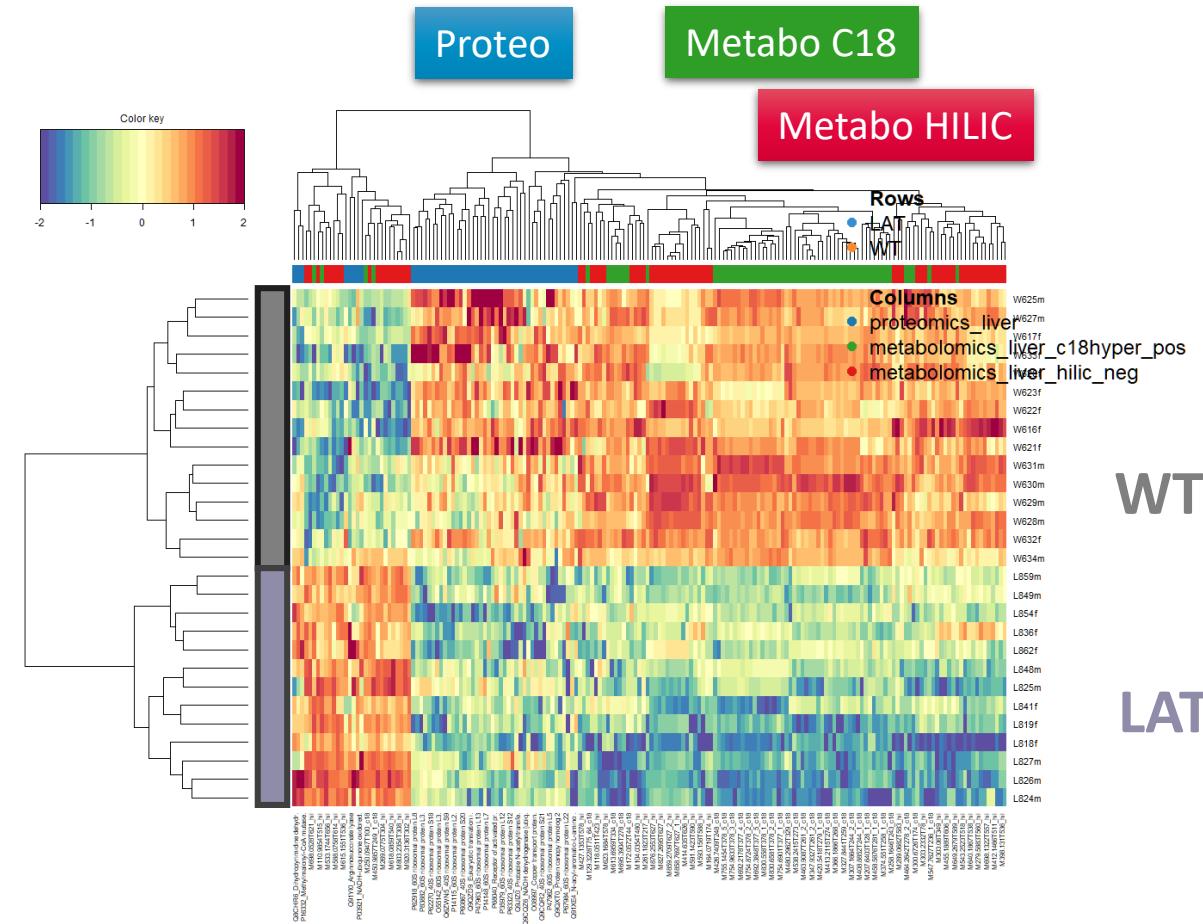


LAT vs WT

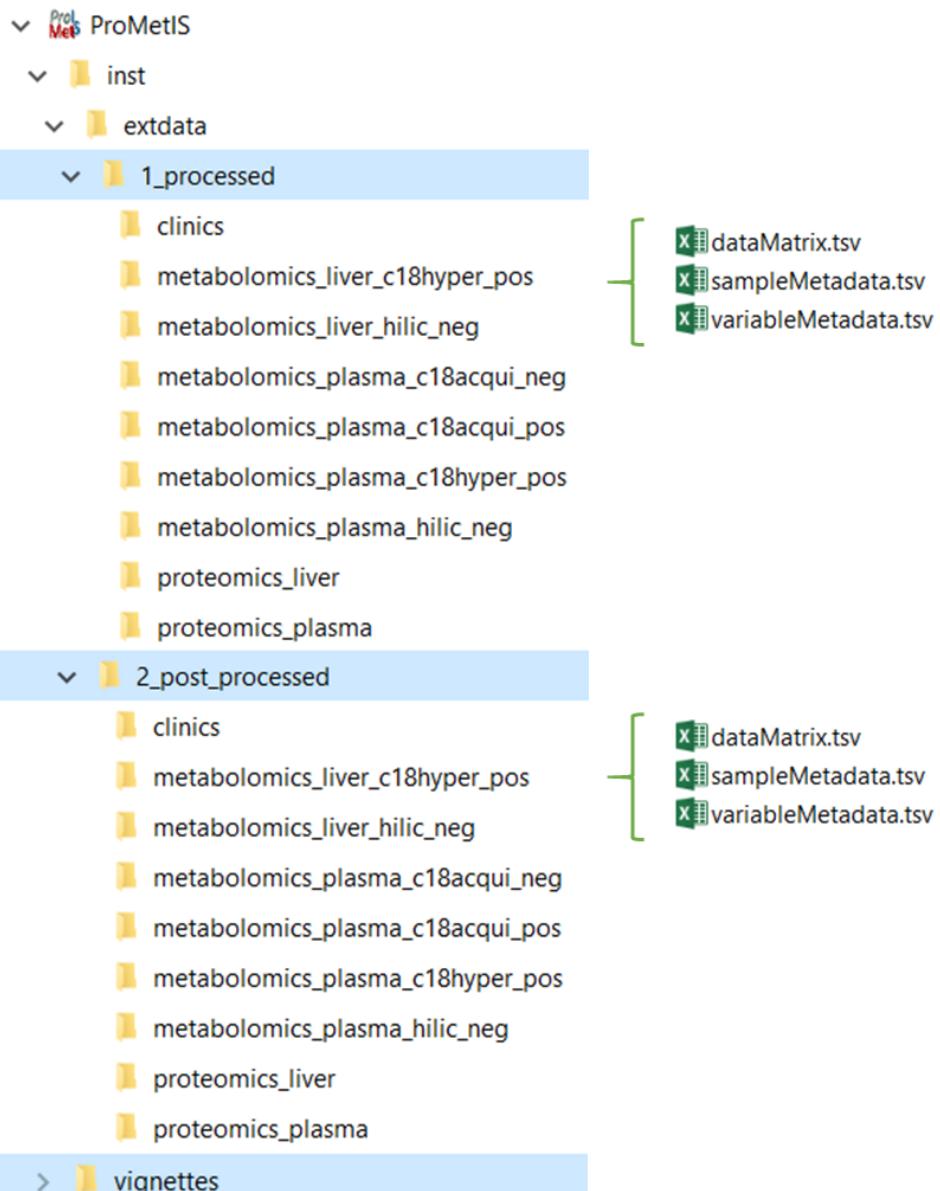


Tenenhaus *et al.* (2014). Variable selection for generalized canonical correlation analysis. *Biostatistics*, **15**:569–583.

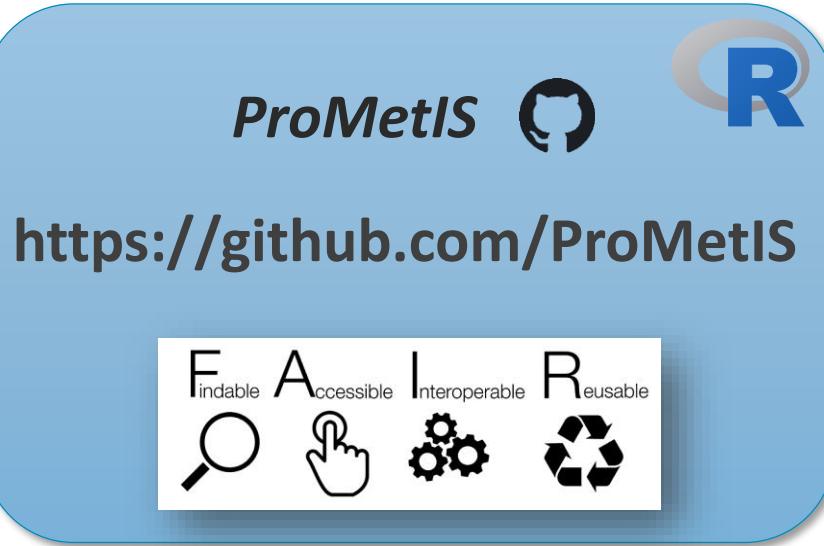
Singh *et al.* (2019). DIABLO: an integrative approach for identifying key molecular drivers from multi-omics assays. *Bioinformatics*, **35**:3055–3062.



Datasets & code availability: the ProMetIS package



Statistical Analysis



3_statistics_singleomics.html
3_statistics_singleomics.Rmd
4_statistics_integrative.html
4_statistics_integrative.Rmd

- ▶ **Value of combining proteomics and metabolomics for fundamental and applied research**
- ▶ **Proteomics and metabolomics data analysis is mature enough to build common pipelines**
- ▶ **Major challenges remain**
 - Limited number of public datasets
 - Limited metabolite annotation
 - Multidisciplinarity

Alyssa Imbert
Camille Roquencourt
Camilo Broc
Krystyna Biletska
Pierrick Roger
Eric Venot
Etienne Thévenot



<https://scidophenia.github.io/>

SciDophenIA



Florence Castelli
Christophe Junot
François Fenaille

Marion Brandolini
Charlotte Joly
Estelle Pujos-Guillot



Magali Rompais
Christine Carapito

Emmanuelle Mouton
Anne Gonzalez de Peredo

Thomas Burger
Yves Vandebrouck
Myriam Ferro



Mohammed Selloum
Tania Sorg
Sophie Leblanc
Yann Herault



Olivier Sand
Jacques van Helden
Claudine Médigue



ProMetIS

SoftwAiR

