



Metabolomics, an introduction with focus on data analysis

International Agency for Research on Cancer
Lyon, France

Reza Salek

International Agency for Research on Cancer

Some Definitions

metabolome

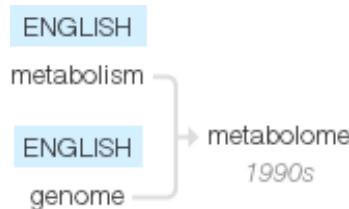
/mə'tabələm/

noun

BIOCHEMISTRY

noun: **metabolome**; plural noun: **metabolomes**

the total number of metabolites present within an organism, cell, or tissue.



metabolomics

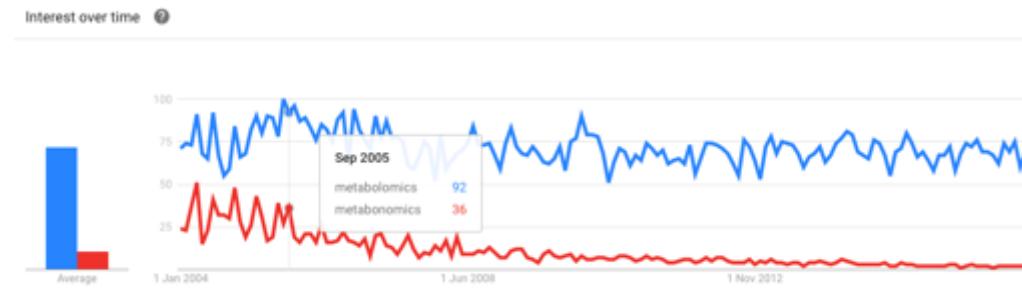
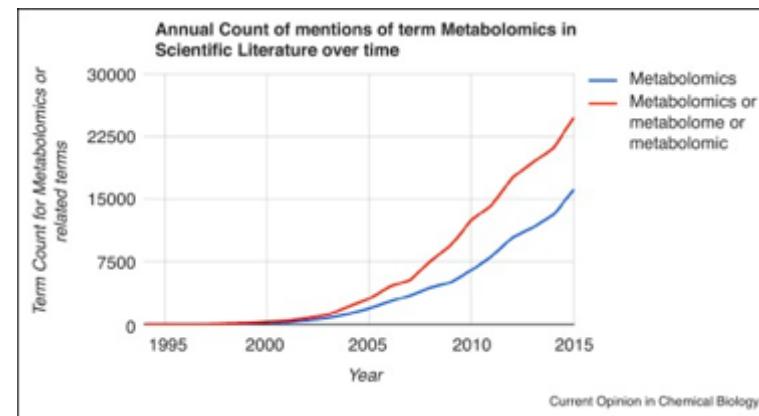
mə'tabələmɪks/

Noun

BIOCHEMISTRY

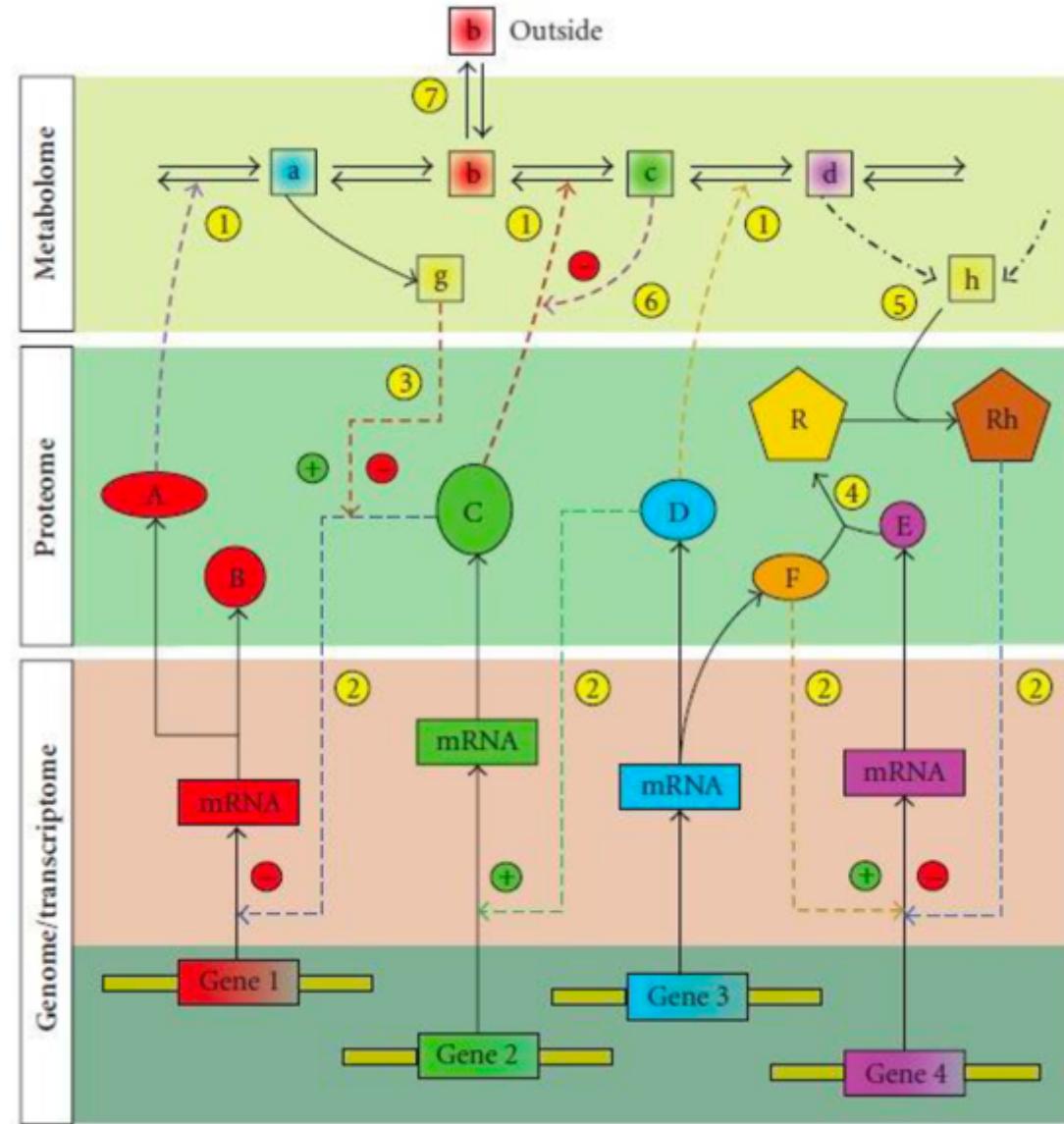
noun: **metabolomics**

the scientific study of the set of metabolites present within an organism, cell, or tissue, often by measuring simultaneously (100s -10000s), many of which are not identified (features or analytes,)

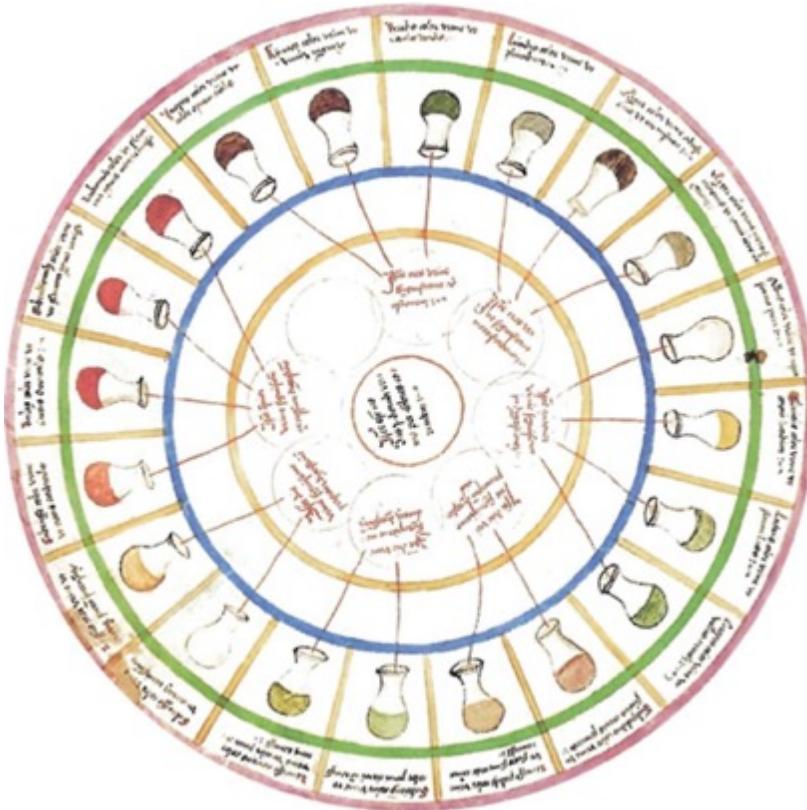


Genomics to metabolomics

- The system (cell) is more than the sum of its parts
- To understand the system we must study the system not the parts
- e.g. Silkworm metamorphosis



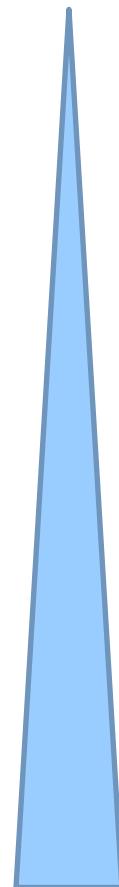
Uroscopy – Early metabolomics



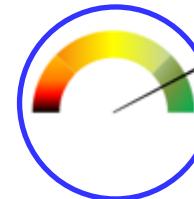
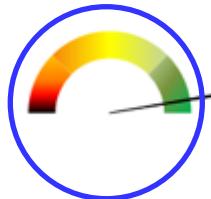
- It dates back to ancient Egypt, Babylon, and India. It was particularly emphasized in Byzantine medicine.
- The wheel describes the possible colors, smells and tastes of urine, and uses them to diagnose disease.

Comparing different metabolomics approaches

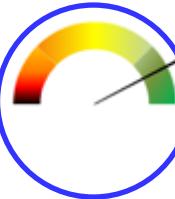
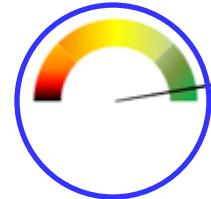
Data
analysis
Complexity



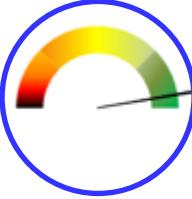
TARGETED ANALYSIS



SUSPECT SCREENING



UNTARGETED SCREENING



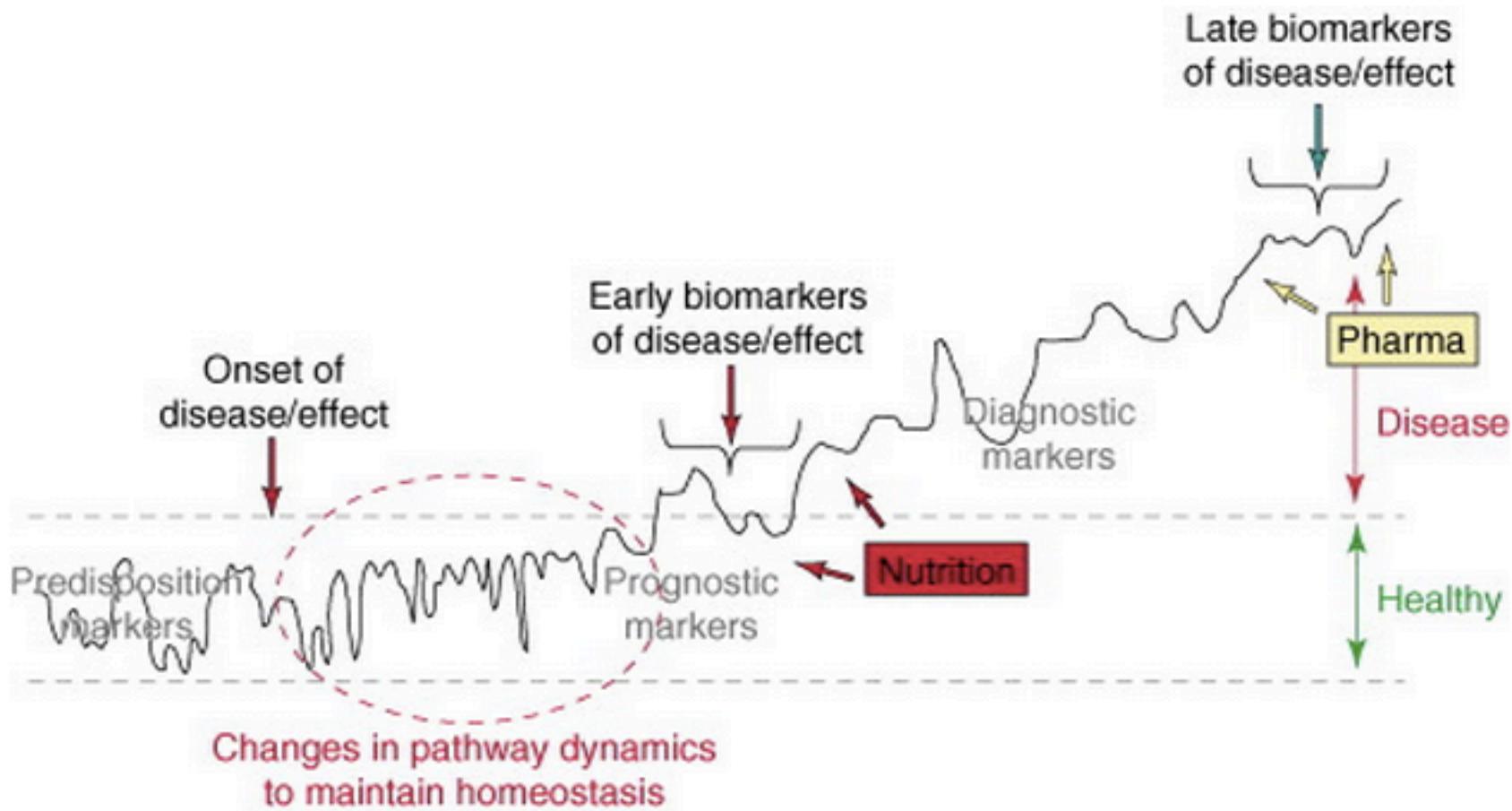
SPEED

SENSITIVITY

COVERAGE

ACCURACY

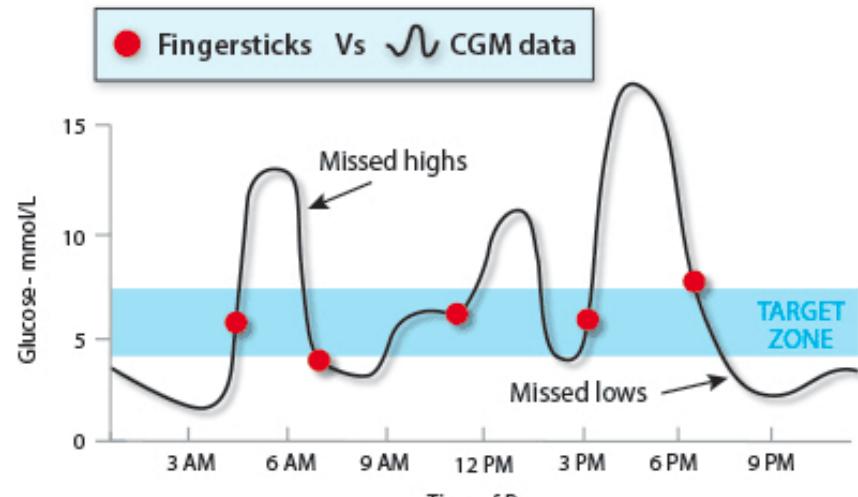
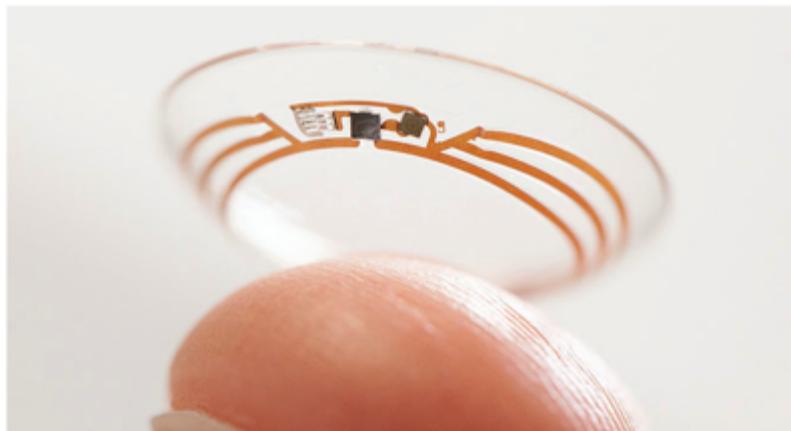
Why is it important



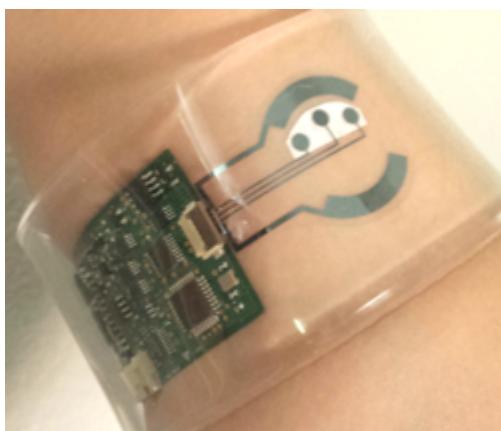
Bio Marker Diagnosis and Monitoring

Google announces 'smart' contact lenses that monitor glucose levels

Published January 16, 2014 / FoxNews.com

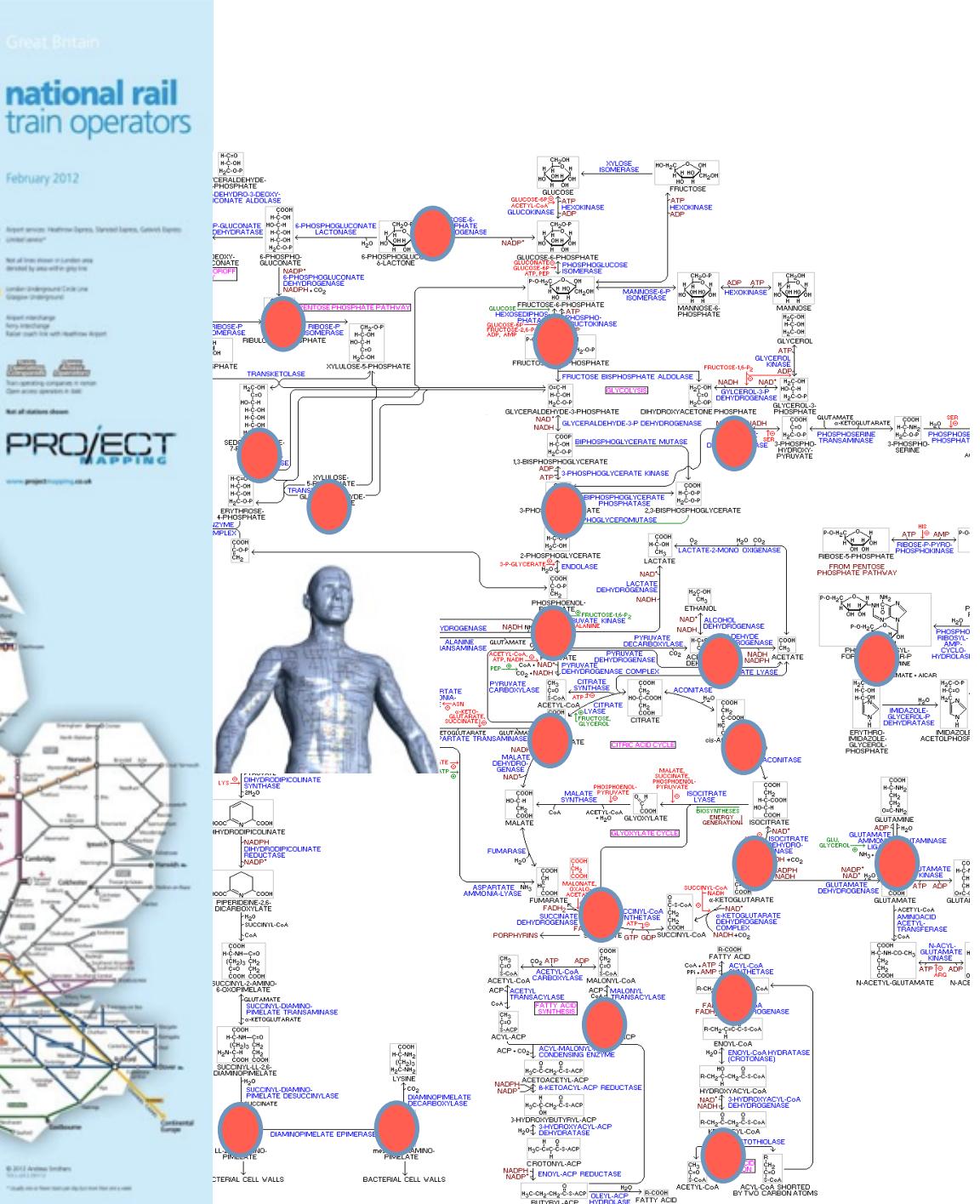


*Illustration purposes only



SMALL WEARABLE SENSOR





Part 1

WHAT WE DO AT IARC

International Agency for Research on Cancer



World Health
Organization

EPIC: European Prospective Investigation into Cancer and nutrition



EPIC Project

- Multicentric cohort study
- Nutrition, lifestyle and cancer
- 23 centres, 10 European countries
- 370.000 women, 150.000 men (aged 35-74) recruited between 1992-1998
- Standardized questionnaires on diet & lifestyle
- Anthropometry measured
- Blood samplers (30ml) on 65% women et 93% men :



8 straws serum
12 straws of plasma
4 straws of buffy coat
4 straws of red blood cells



EWAS and dietary intake - Animal foods

Dietary intervention study

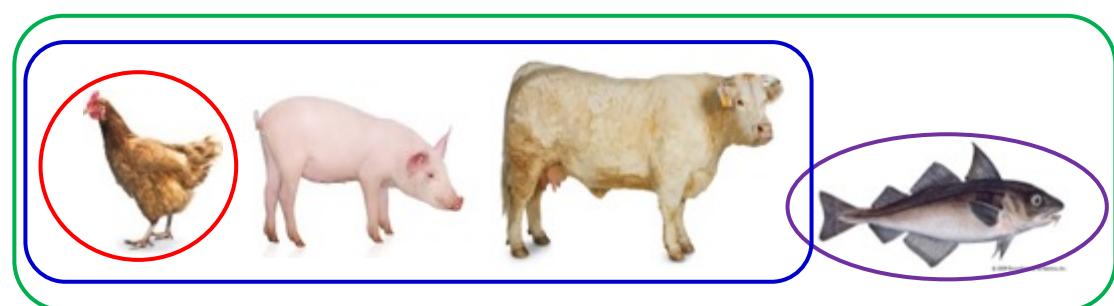
4 parallel groups (n=10 per group)

Controlled diet

3 weeks with increasing level
of meat/fish intake

24-hr urine, serum

UHPLC-QTof-MS



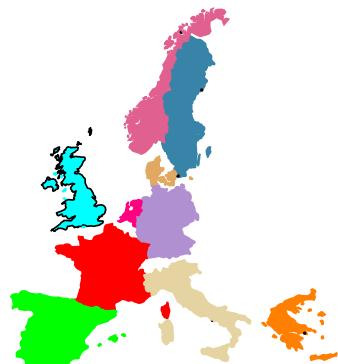
2-Methylbutyrylcarnitine
Acetylcarnitine
Propionylcarnitine

Observational study

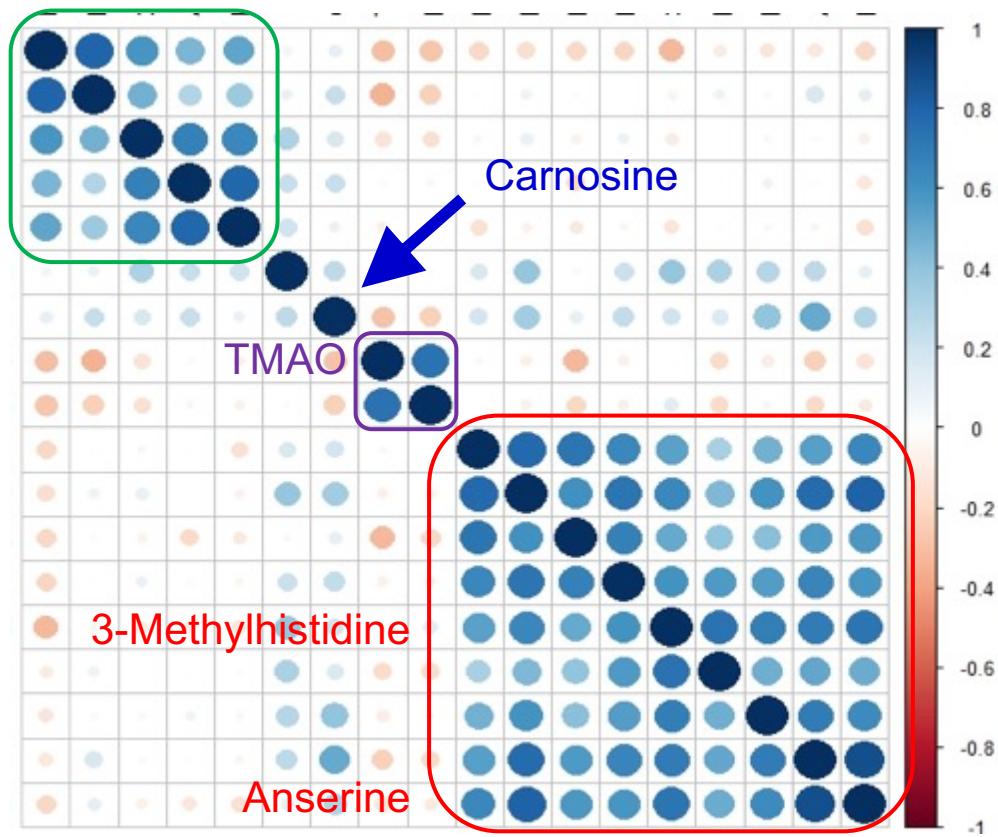
24-hr dietary recalls

5 groups (n=10) according to meat/fish
intake

24-hr urine



International Agency for Research on Cancer



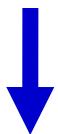
Cheung et al., 2017, Am. J. Clin. Nutr.



Untargeted metabolomics and risk of hepatocellular carcinoma

129 HCC cases
129 matched controls
Plasma samples
Untargeted metabolomics
(LC-QTof-MS)

9,602 MS features
detected



92 discriminant
metabolites



46 metabolites
annotated

Top 16 HCC metabolites	LC-MS Method	Multivariable OR (95% CI)
Retinol	RP+	0.27 (0.16 - 0.48)
Dehydroepiandrosterone sulfate	HILIC-	0.35 (0.22 - 0.57)
Glycerophosphocholine	RP+	0.44 (0.28 - 0.71)
γ -carboxyethylhydroxchroman	RP+	0.56 (0.39 - 0.81)
Creatine	RP+	0.56 (0.37 - 0.83)
Tyrosine	RP+	2.04 (1.30 - 3.20)
N1-Acetylspermidine	HILIC+	2.16 (1.38 - 3.37)
Isatin	RP+	2.28 (1.38 - 3.75)
<i>p</i> -Hydroxyphenyllactic acid	HILIC-	2.77 (1.58 - 4.83)
Sphingosine	RP+	2.79 (1.66 - 4.71)
L,L-Cyclo(leucylprolyl)	RP+	3.25 (1.91 - 5.53)
Glycochenodeoxycholic acid	RP+	3.31 (1.99 - 5.51)
Glycocholic acid	RP+	4.07 (2.32 - 7.14)
7-Methylguanine	HILIC+	6.78 (3.24 - 14.18)

- Many metabolites replicated in in the ATBC cohort (Finland)



Stepien et al., to be submitted



International Agency for Research on Cancer



<http://exposome-explorer.iarc.fr/>



compounds, foods, cancers



[Home](#) Biomarker data ▾ Classifications ▾ Publications Structure search Downloads About ▾



Exposome-Explorer

First database dedicated to biomarkers of exposure to environmental risk factors for diseases.

[Read more](#)



Neveu V, Nicolas G, Salek RM, Wishart DS, Scalbert A. Exposome-Explorer 2.0., 2019, *Nucl. Acid Res.*

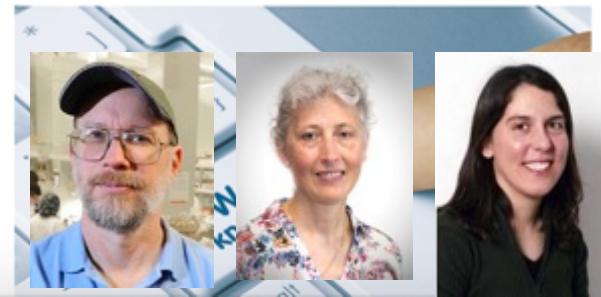
Biomarker data



Classifications



Additional information



Raw Data to Peak tables

METABOLOMICS DATA

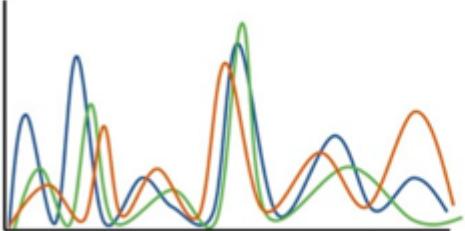
International Agency for Research on Cancer



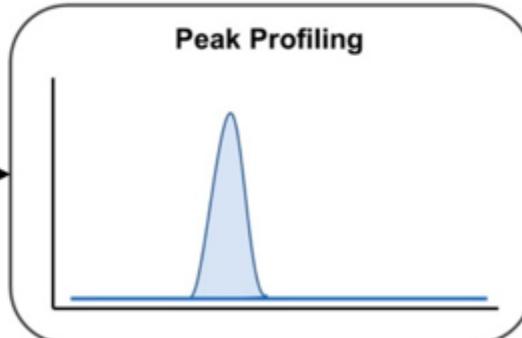
Overview of data handling in metabolomics

Raw Data Processing

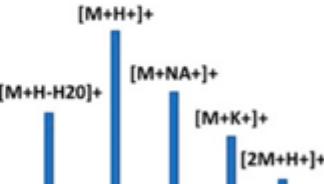
LC-MS Spectra



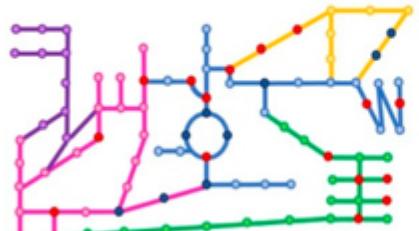
Peak Profiling



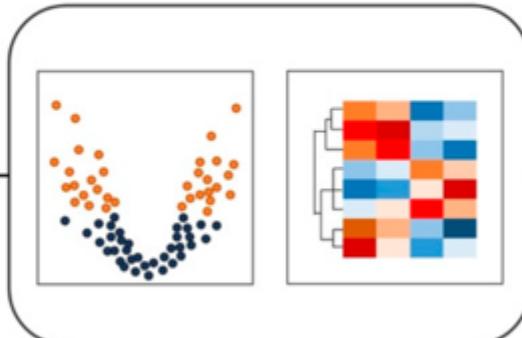
Peak Annotation



Functional Analysis



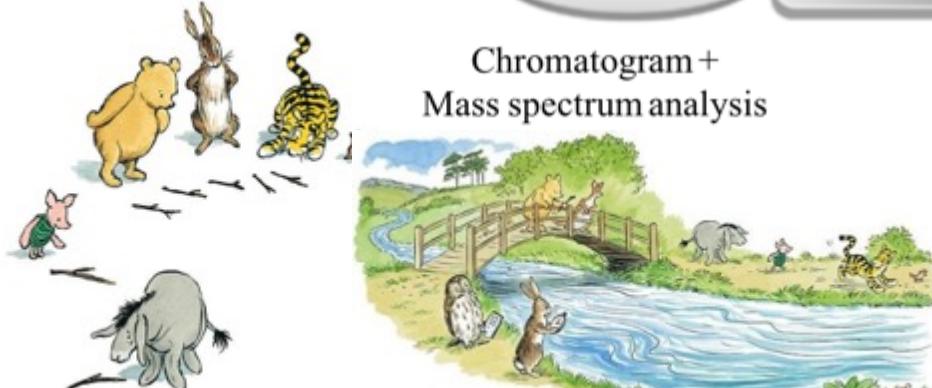
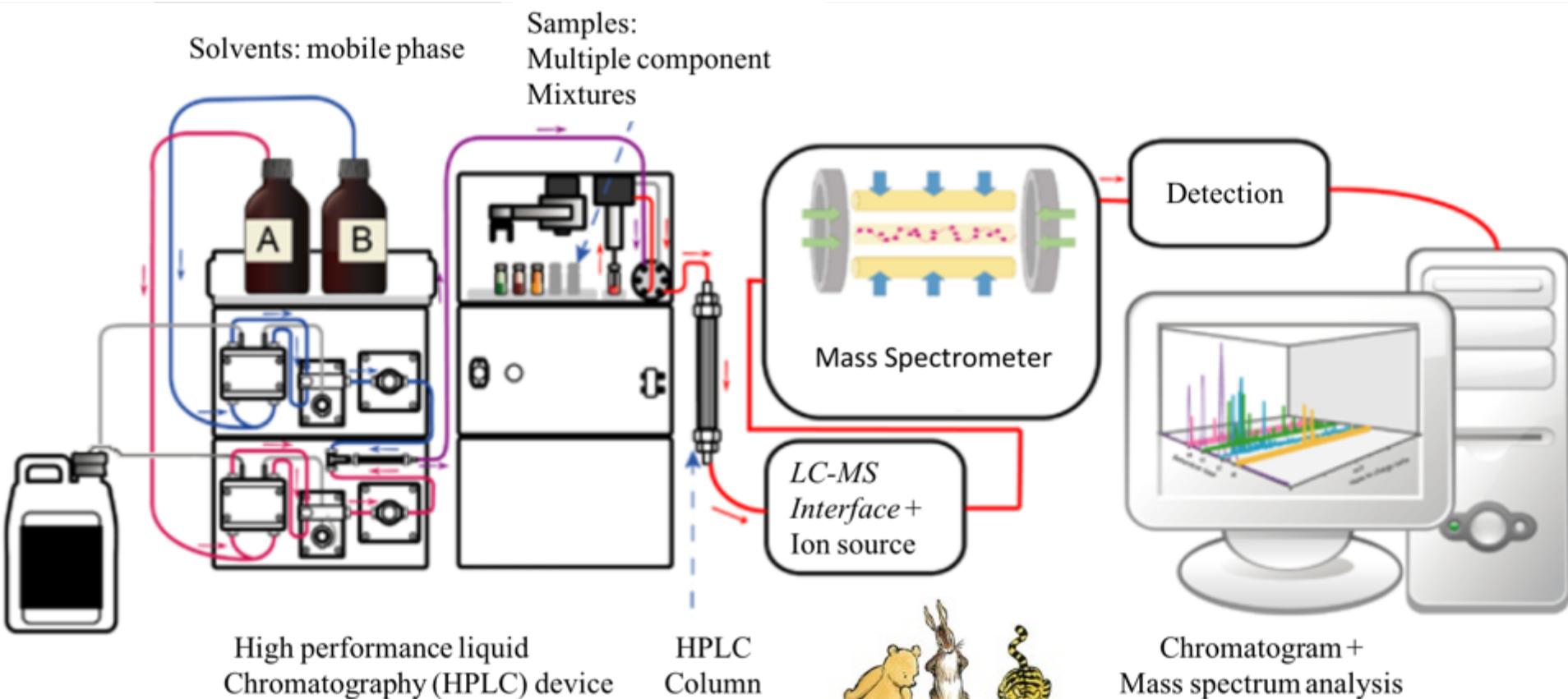
Statistical Analysis



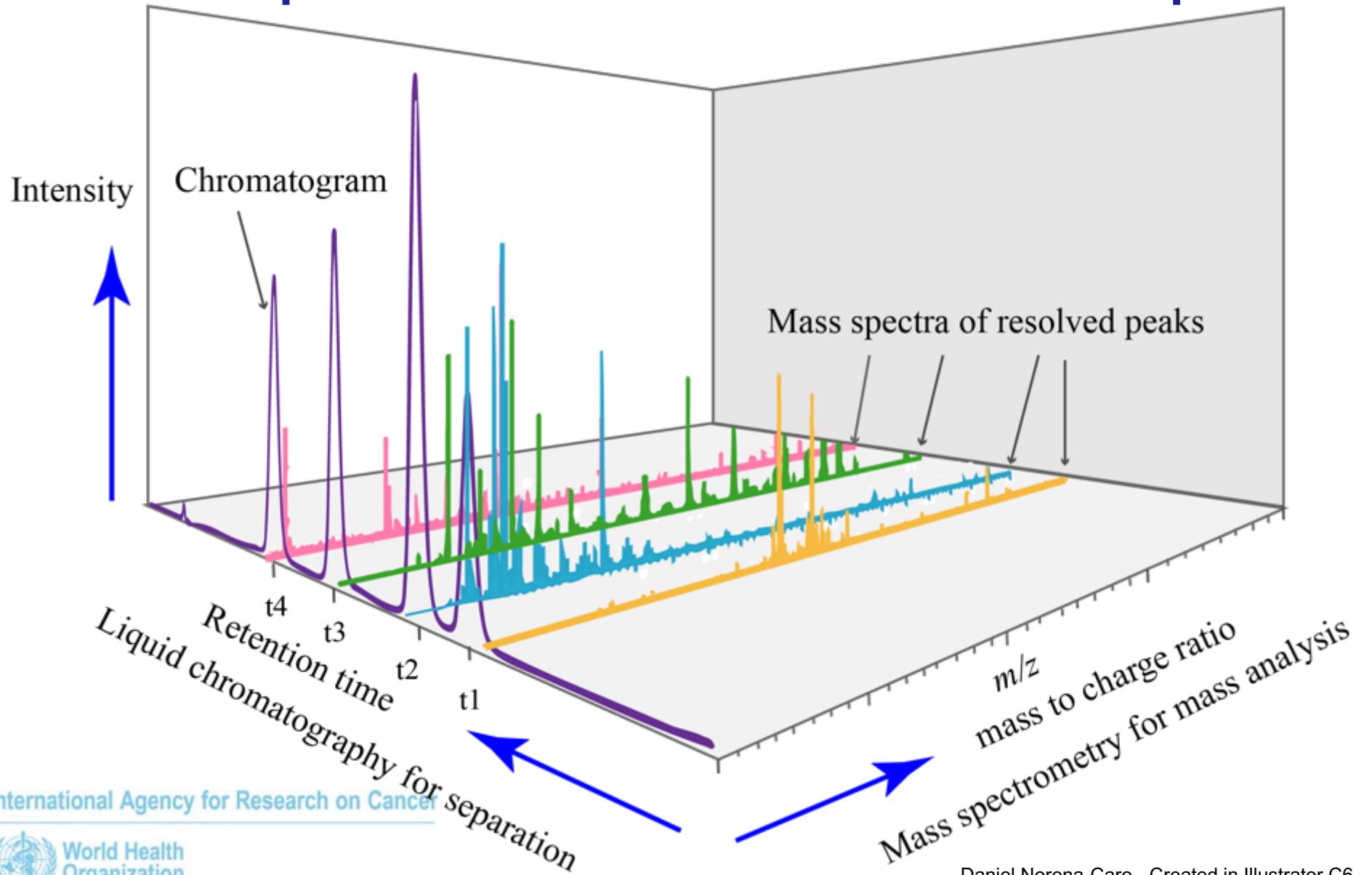
Data Cleaning

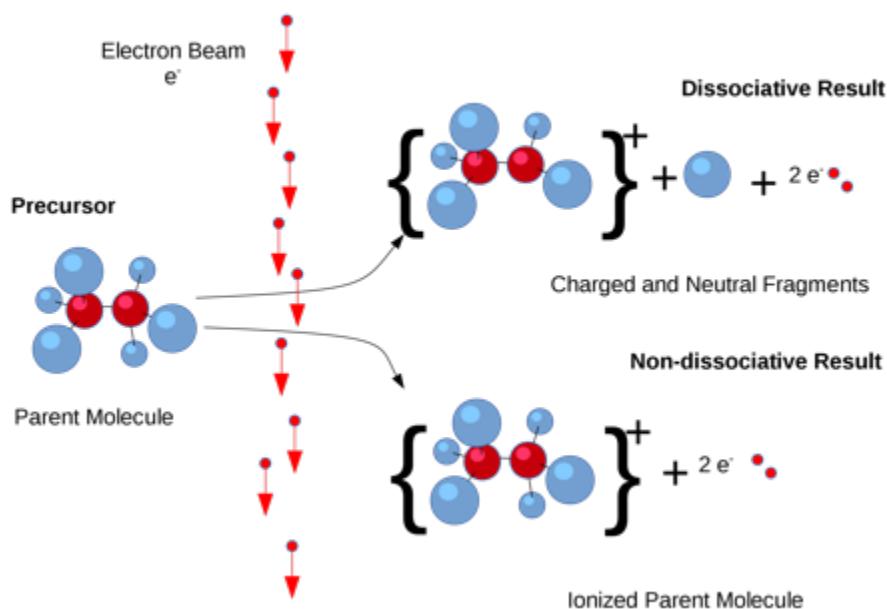
m/z	Sample A	Sample B
99.0345	10145.22	2825.35
108.7492	28491.08	5398.11
213.8490	531.89	15410.92

Diagram of an LC-MS system



LC-MS spectrum of each resolved peak

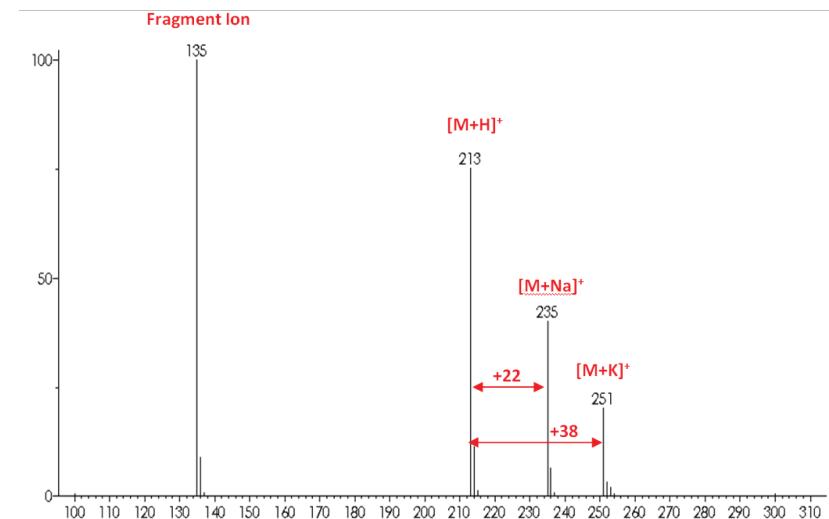




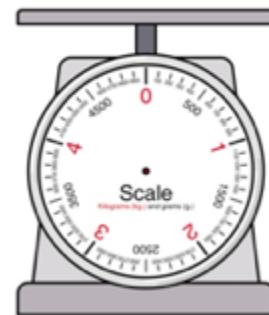
$$H1.0078 \text{ u} - e 0.00054 \text{ u} = 1.0073$$

Ionization	Formation	Ion Mass
Positive	$[M+H]^+$	$m+1.0073$
	$[M+2H]^{2+}$	$m/2+1.0073$
	$[M+Na]^+$	$m+22.9892$
	$[M+K]^+$	$m+38.9632$
	$[M+NH_4]^+$	$m+18.03382$
Negative	$[M-H]^-$	$m-1.0073$
	$[M-2H]^{2-}$	$m/2-1.0073$
	$[M-2H+Na]^-$	$m+20.9747$
	$[M-2H+K]^-$	$m+36.9486$

M is the molecule with molecular weight m

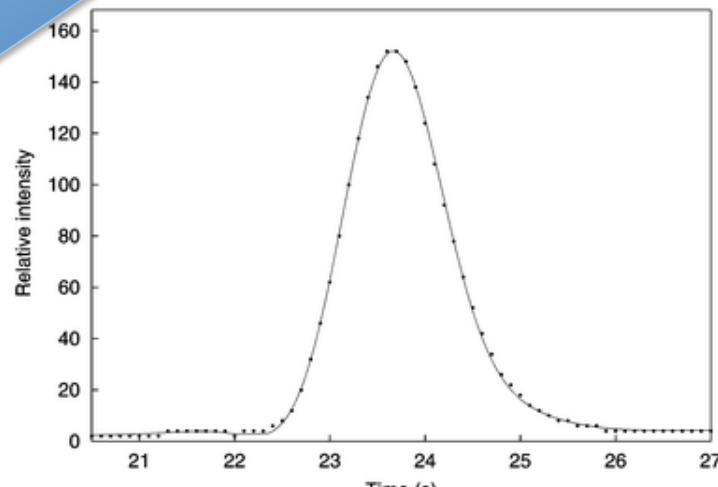
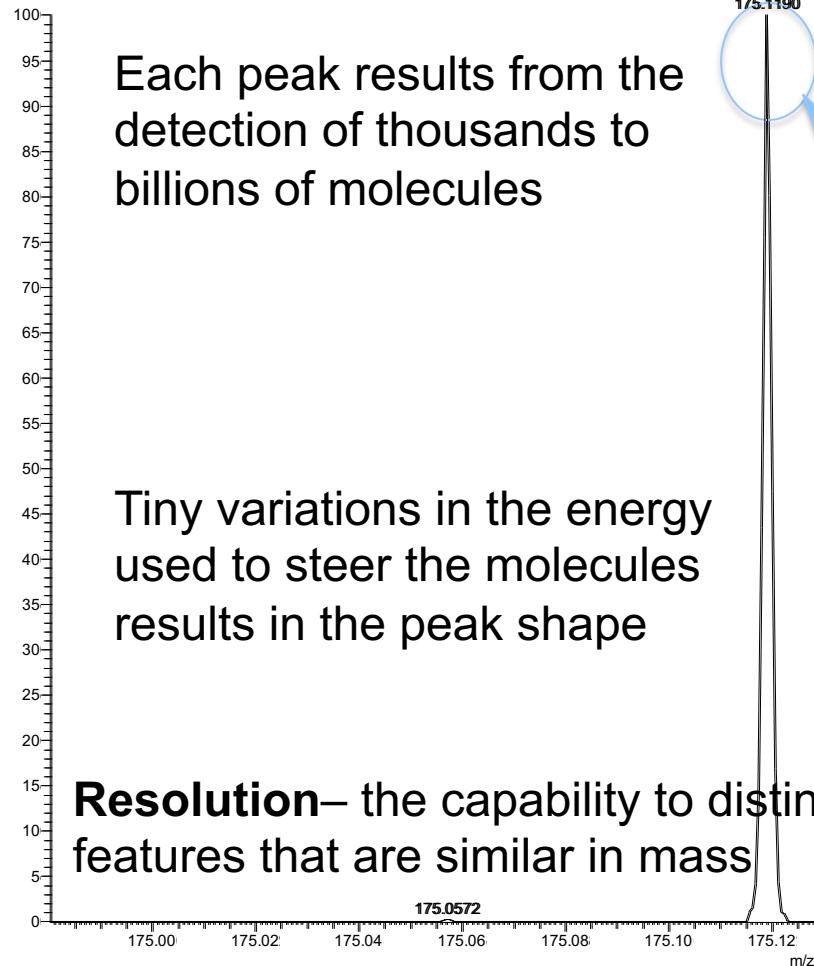


Molecules and fragments

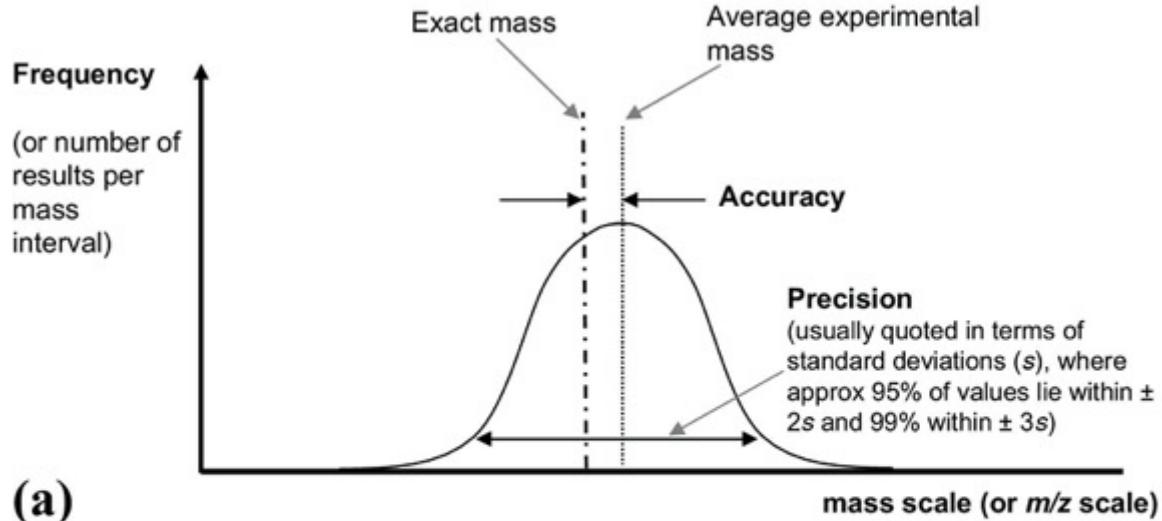


MS Spectra

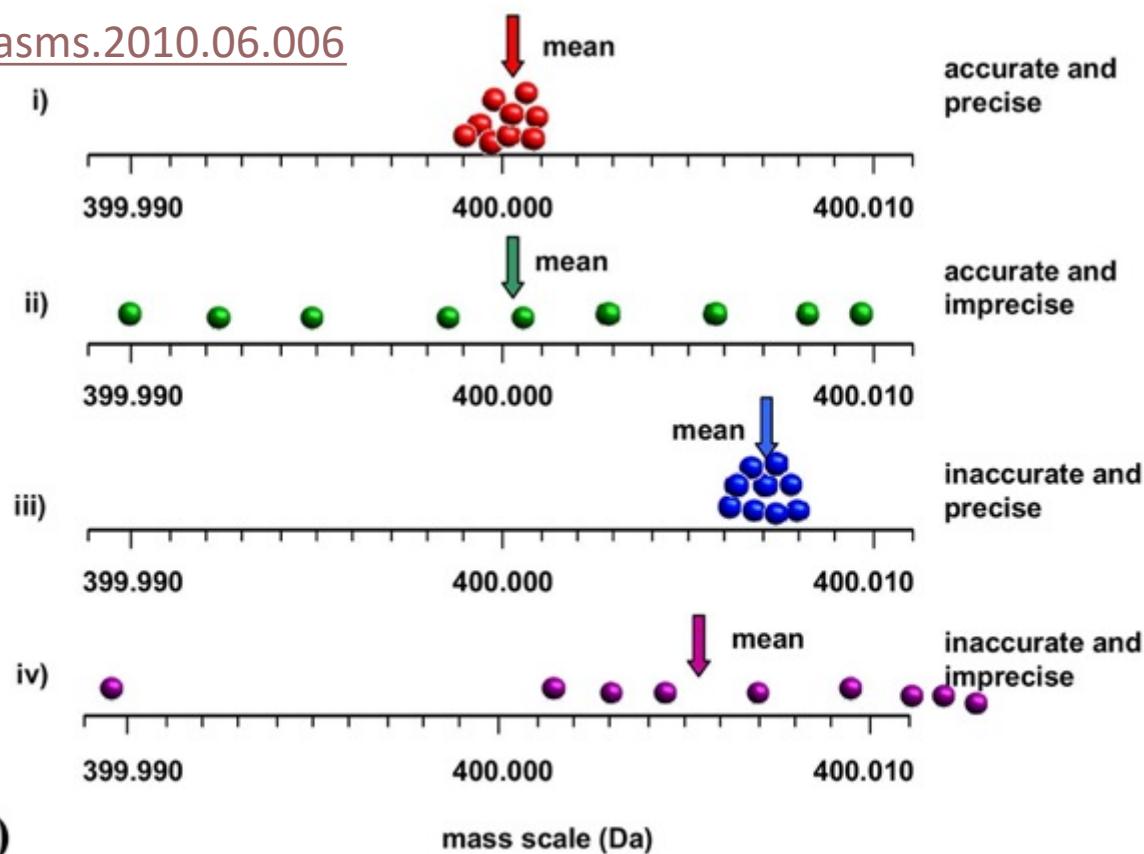
QCB_16FEB18_2 #1827 RT: 19.50 AV: 1 NL: 9.12E6
T: FTMS + p ESI Full lock ms [70.0000-1050.0000]



Each point in this peak graph is an MS scan



<https://doi.org/10.1016/j.jasms.2010.06.006>



Exploiting High Mass Accuracy to ID Compounds

<u>Type</u>	<u>Mass Accuracy</u>
FT-ICR-MS	0.1 - 1 ppm
Orbitrap	0.5 - 1 ppm
Magnetic Sector	1 - 2 ppm
TOF-MS	3 - 5 ppm
Q-TOF	3 - 5 ppm
Triple Quad	3 - 5 ppm
Linear IonTrap	50-200 ppm (10 ppm in Ultra-Zoom)

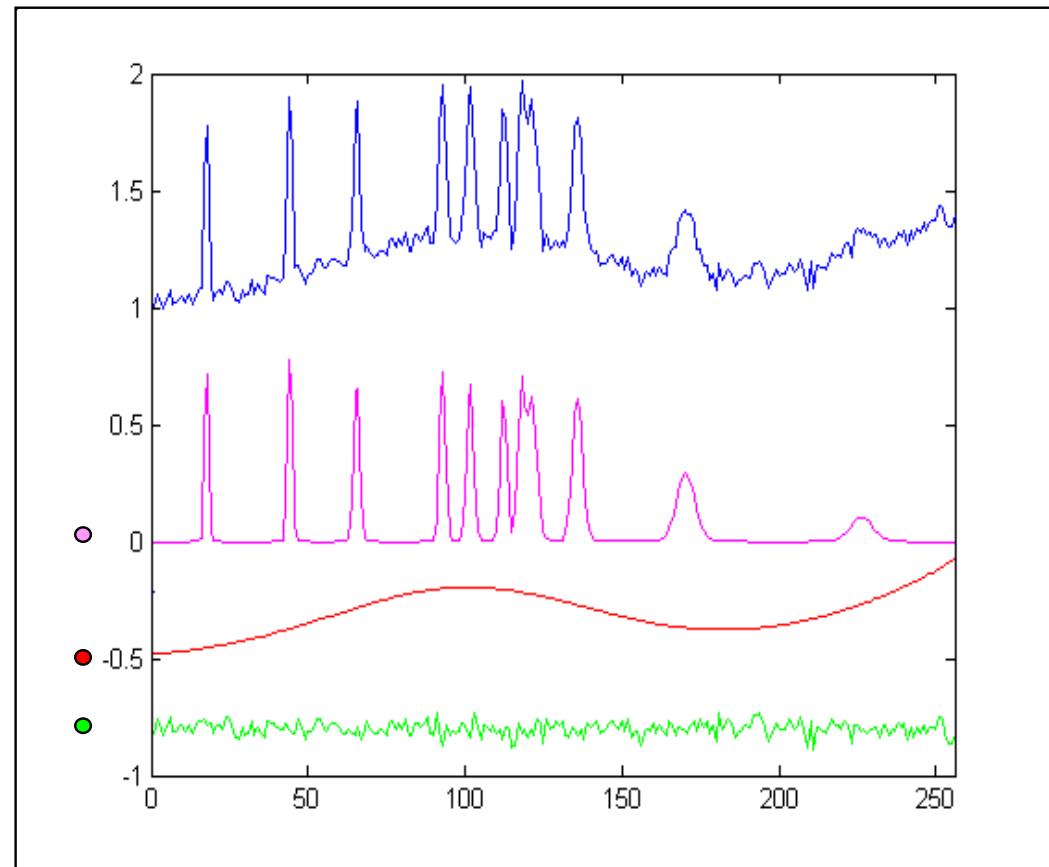
$$\text{ppm} = \left(\frac{\text{m}_{\text{exp}} - \text{m}_{\text{calc}}}{\text{m}_{\text{exp}}} \right) * 1 \text{E} + 6$$

Signal/Spectra components

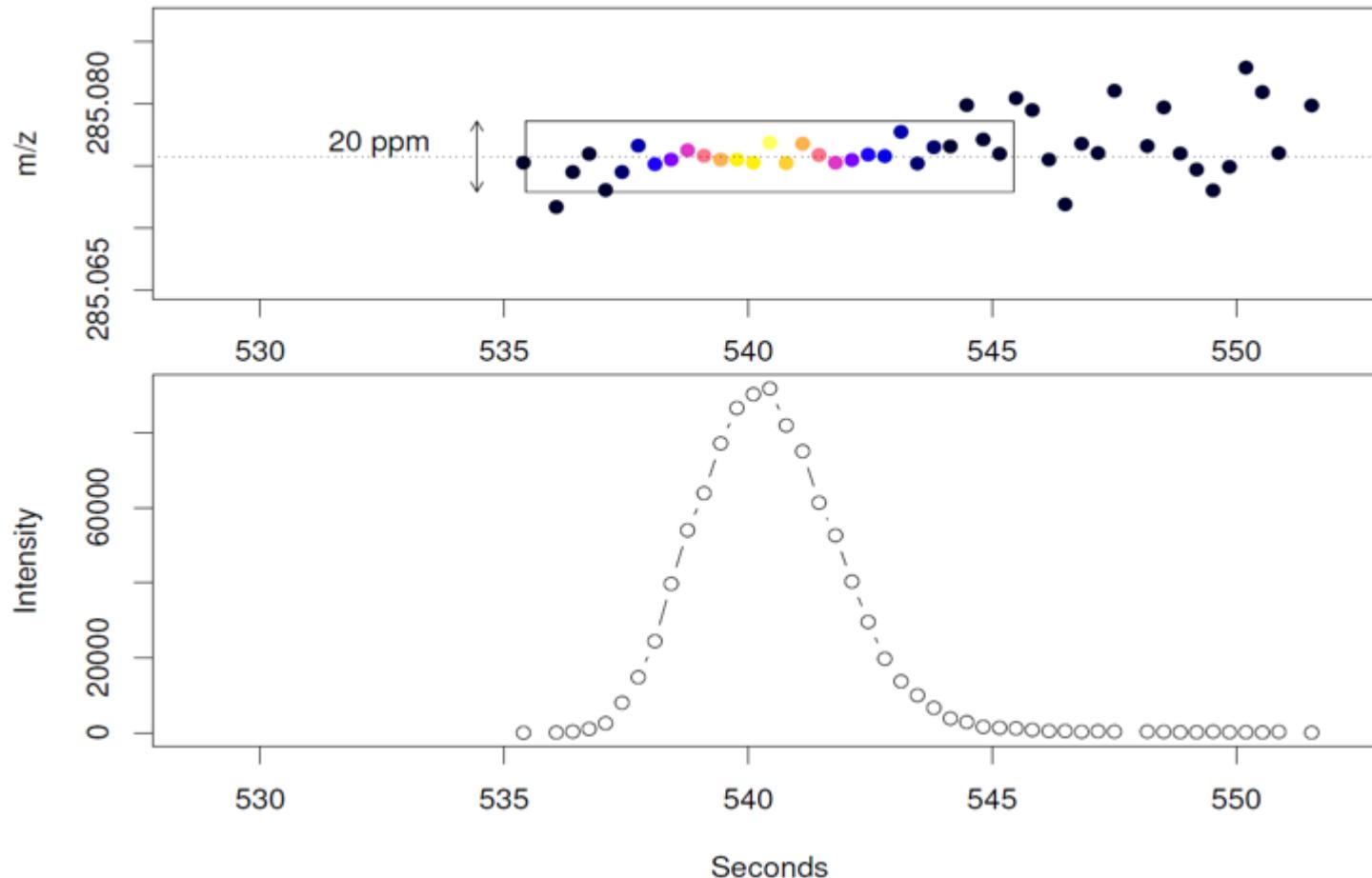
- signal
- noise
- background

Signal enhancement

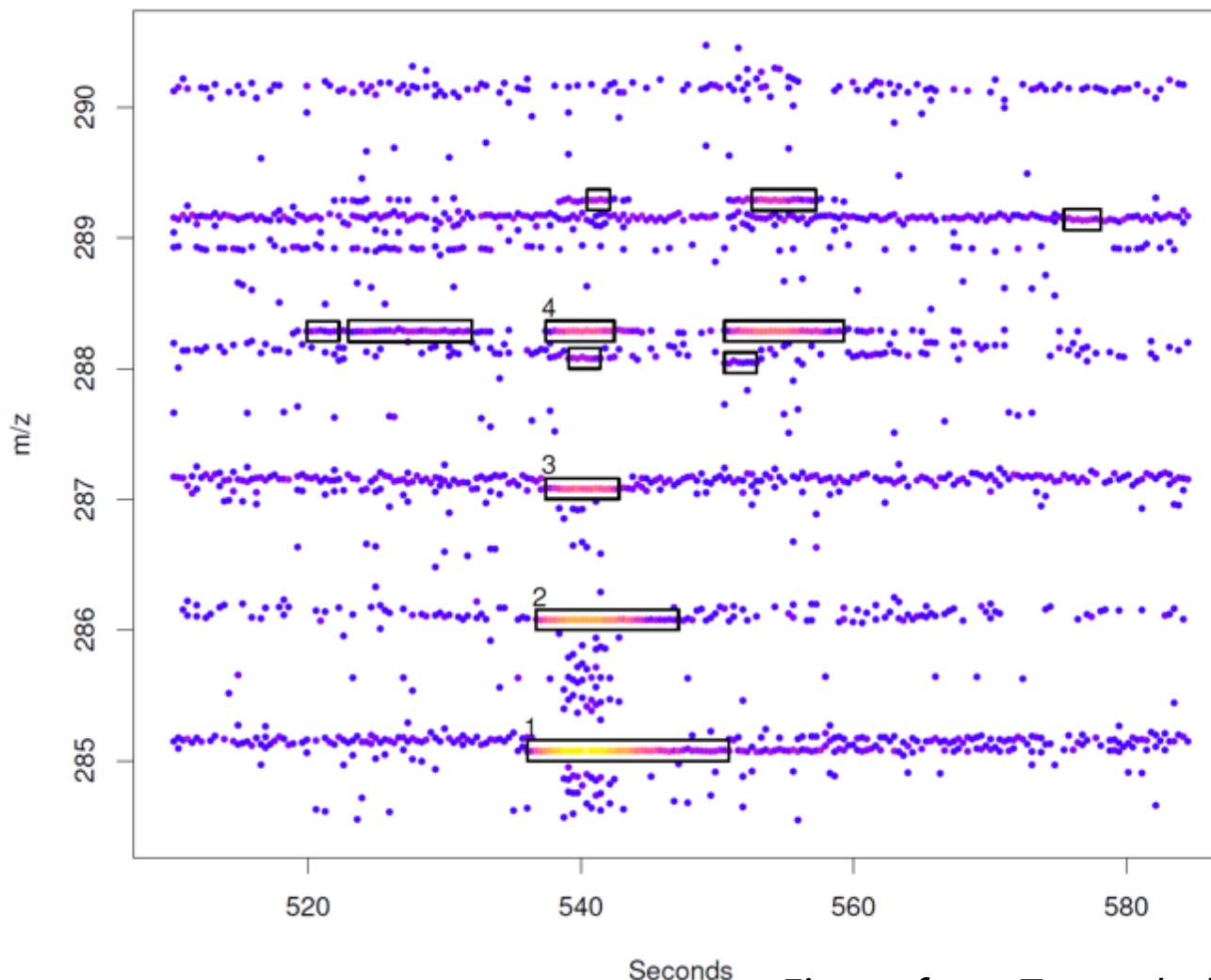
De-noising
Background correction



Example: XCMS Centwave algorithm

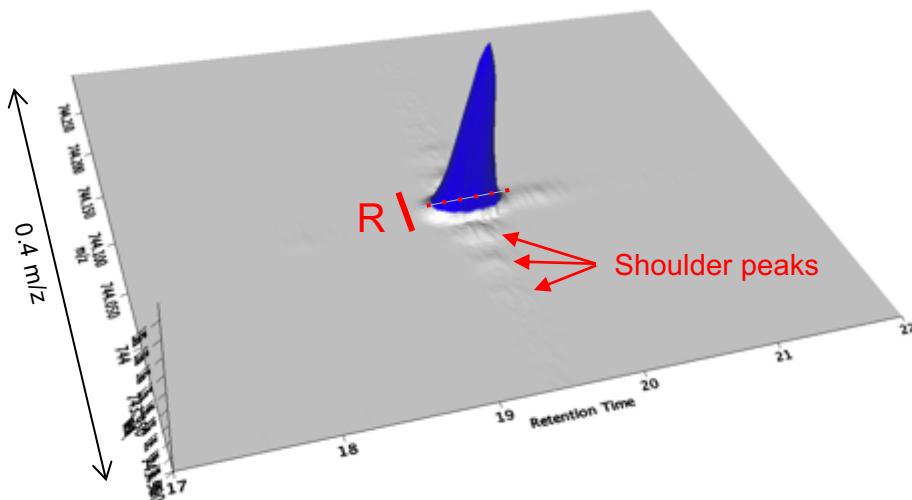


XCMS centwave algorithm

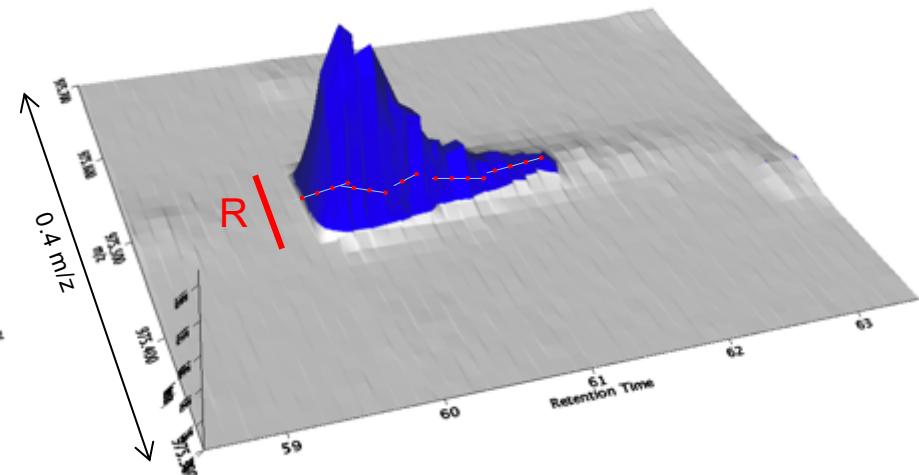


A single ion seen on different MS instruments

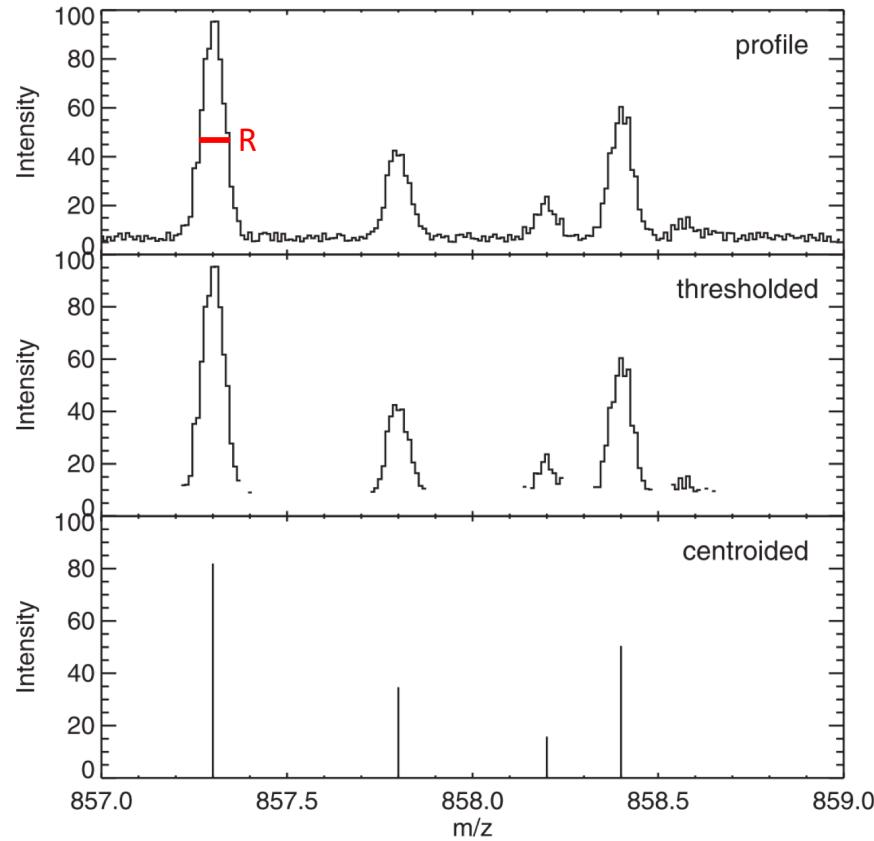
Thermo Orbitrap



Waters SYNAPT G2-S



Profile mode vs Centroid



Deutsch, *Mol. Cell. Proteomics* 2012

International Agency for Research on Cancer

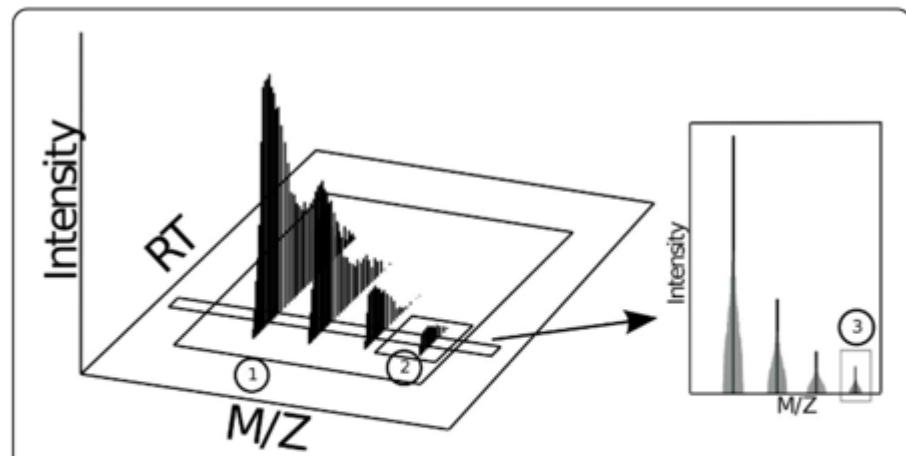
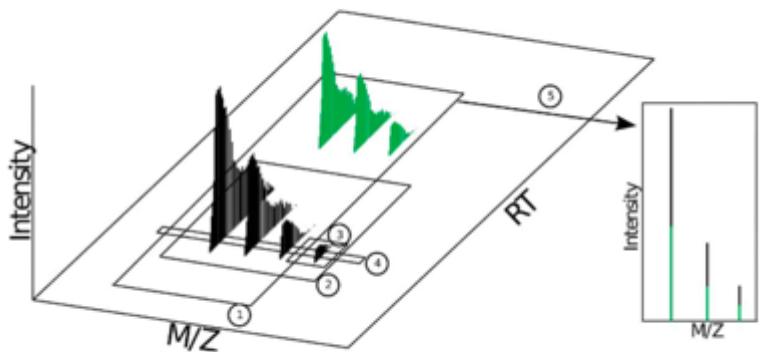


Tomáš Pluskal

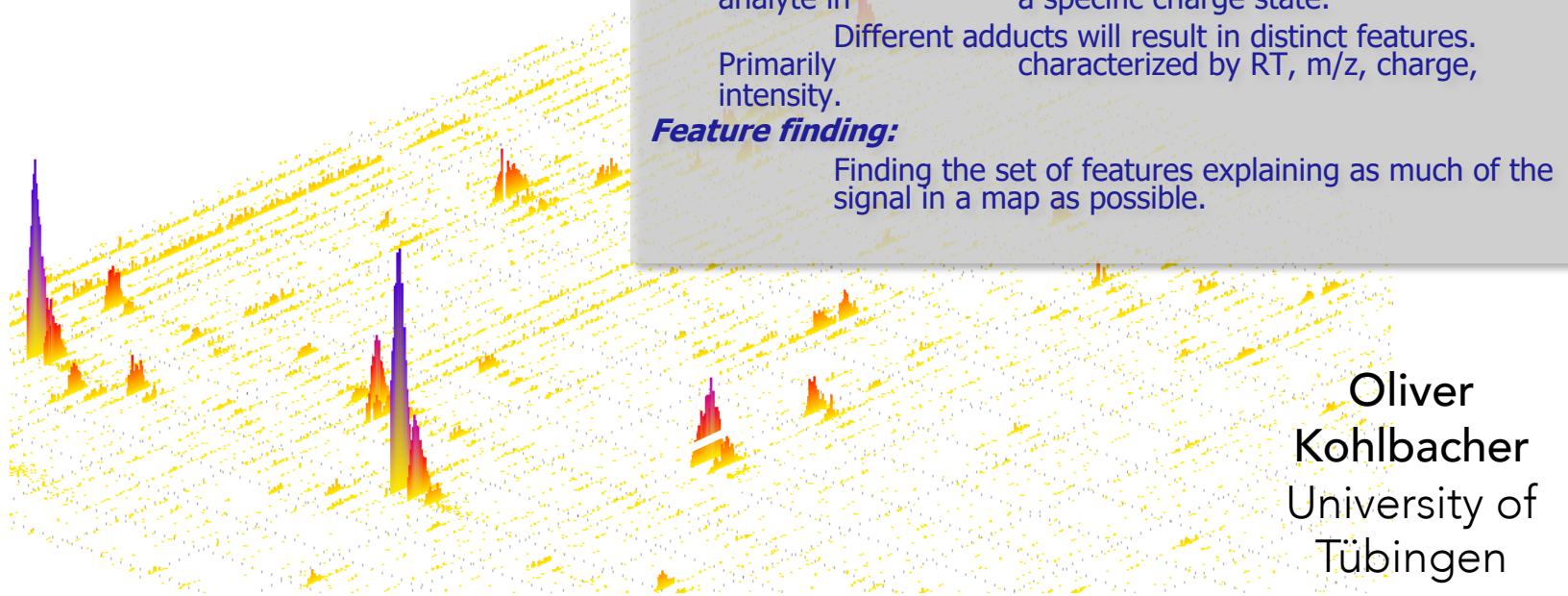
Terminology artifacts

Feature
Peak
Component
Isotope trace
Isotope envelope

Picking
Finding
Detection
Extraction
Deconvolution

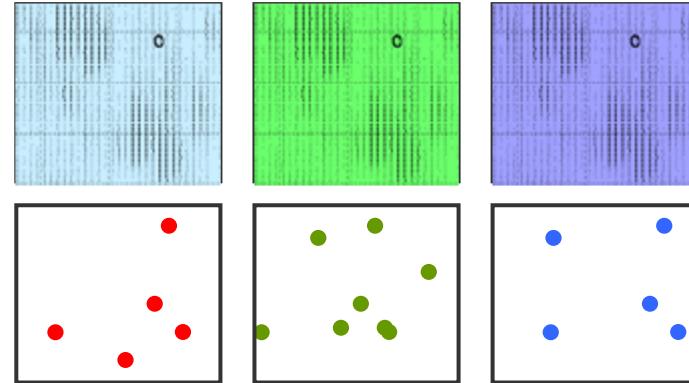


Feature Finding – Terms



Metabolic Profiling

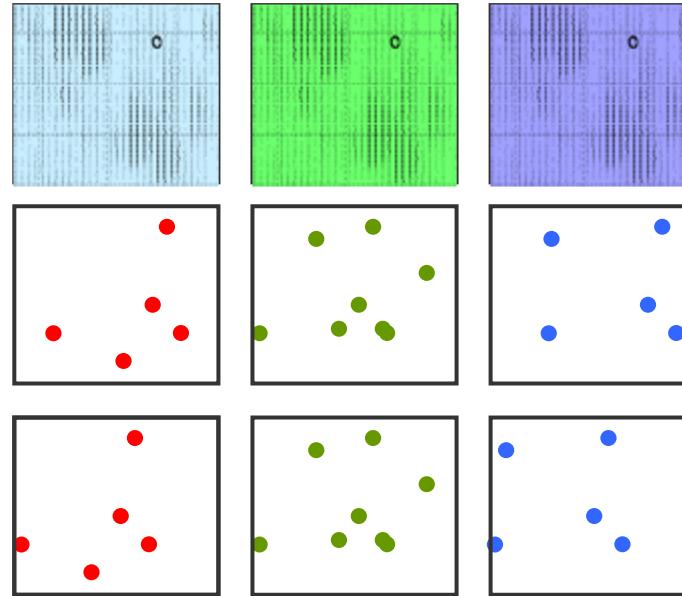
1. Find features in all maps



Oliver
Kohlbacher
University of
Tübingen

Metabolic Profiling

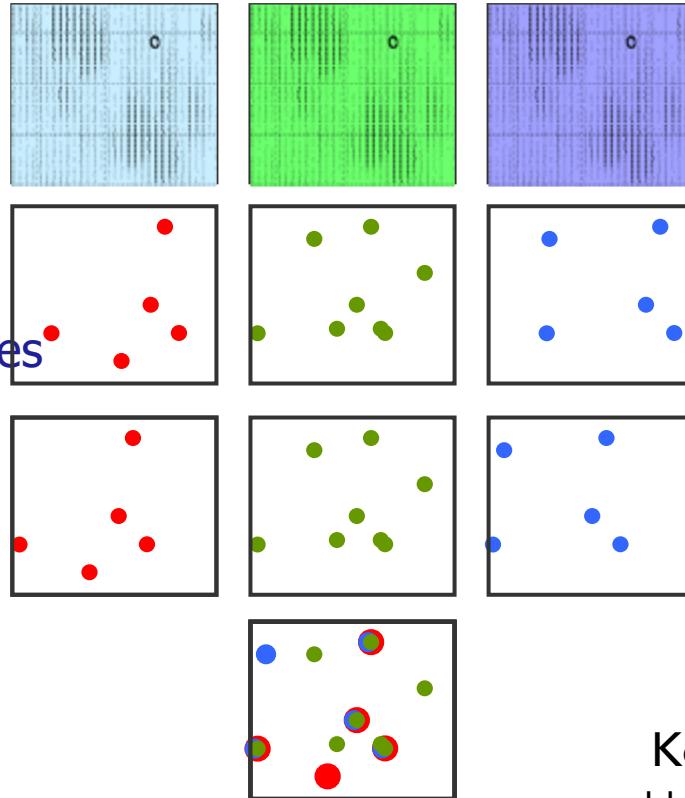
1. **Find** features in all maps
2. **Align** maps



Oliver
Kohlbacher
University of
Tübingen

Metabolic Profiling

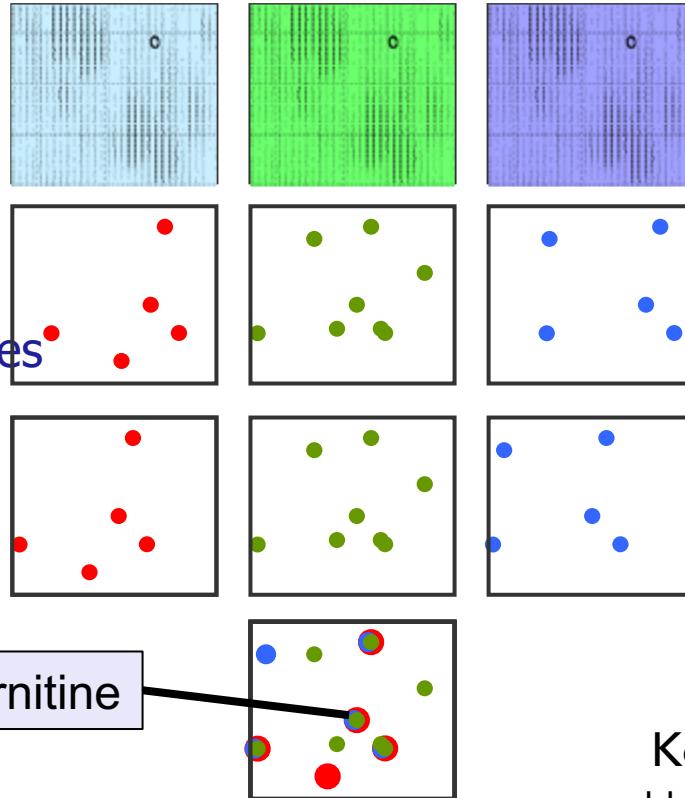
1. **Find** features in all maps
2. **Align** maps
3. **Link** corresponding features



Oliver
Kohlbacher
University of
Tübingen

Metabolic Profiling

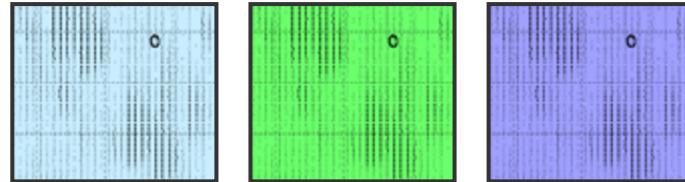
1. **Find** features in all maps
2. **Align** maps
3. **Link** corresponding features
4. **Identify** features



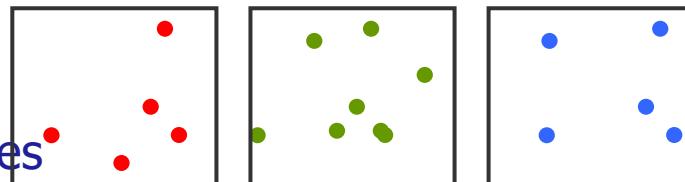
Oliver
Kohlbacher
University of
Tübingen

Metabolic Profiling

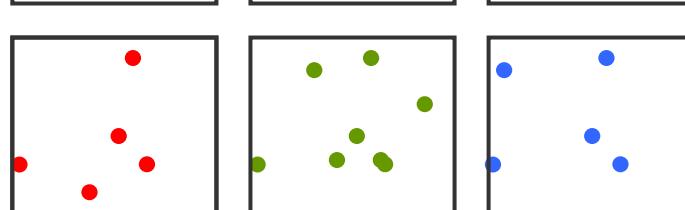
1. **Find** features in all maps



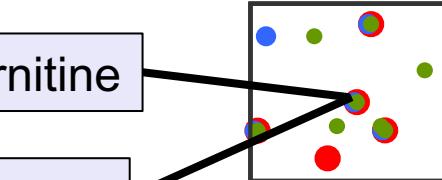
2. **Align** maps



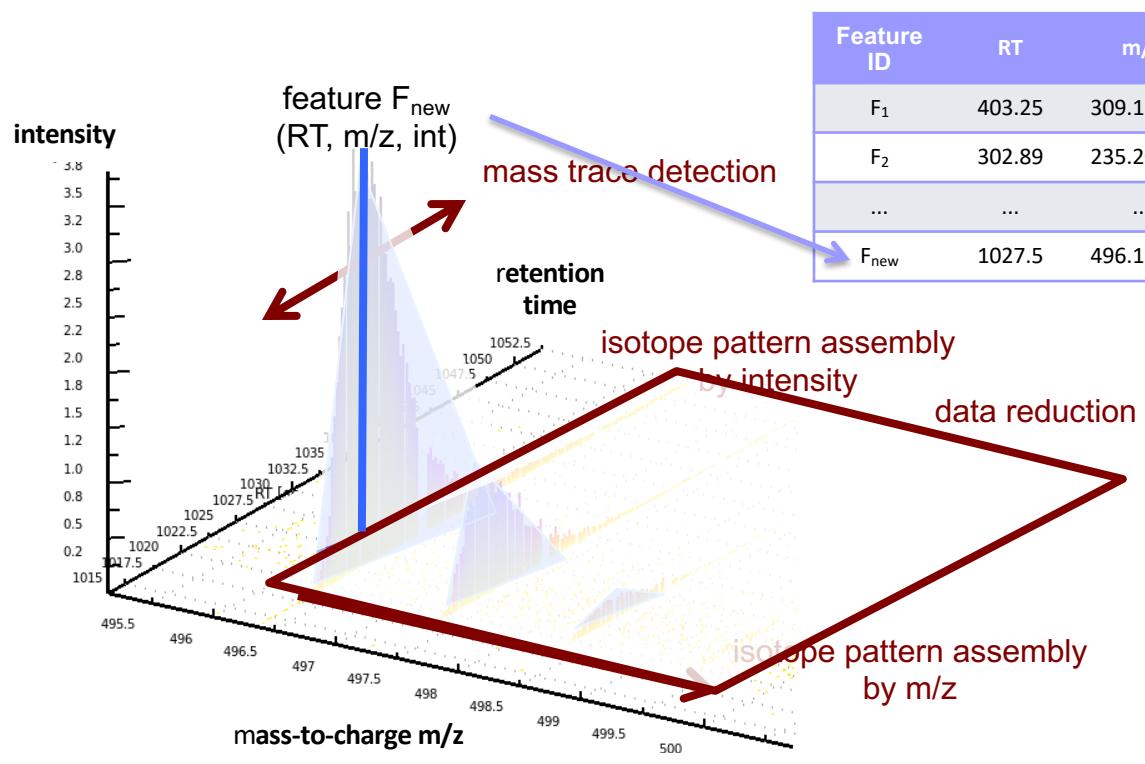
3. **Link** corresponding features



4. **Identify** features



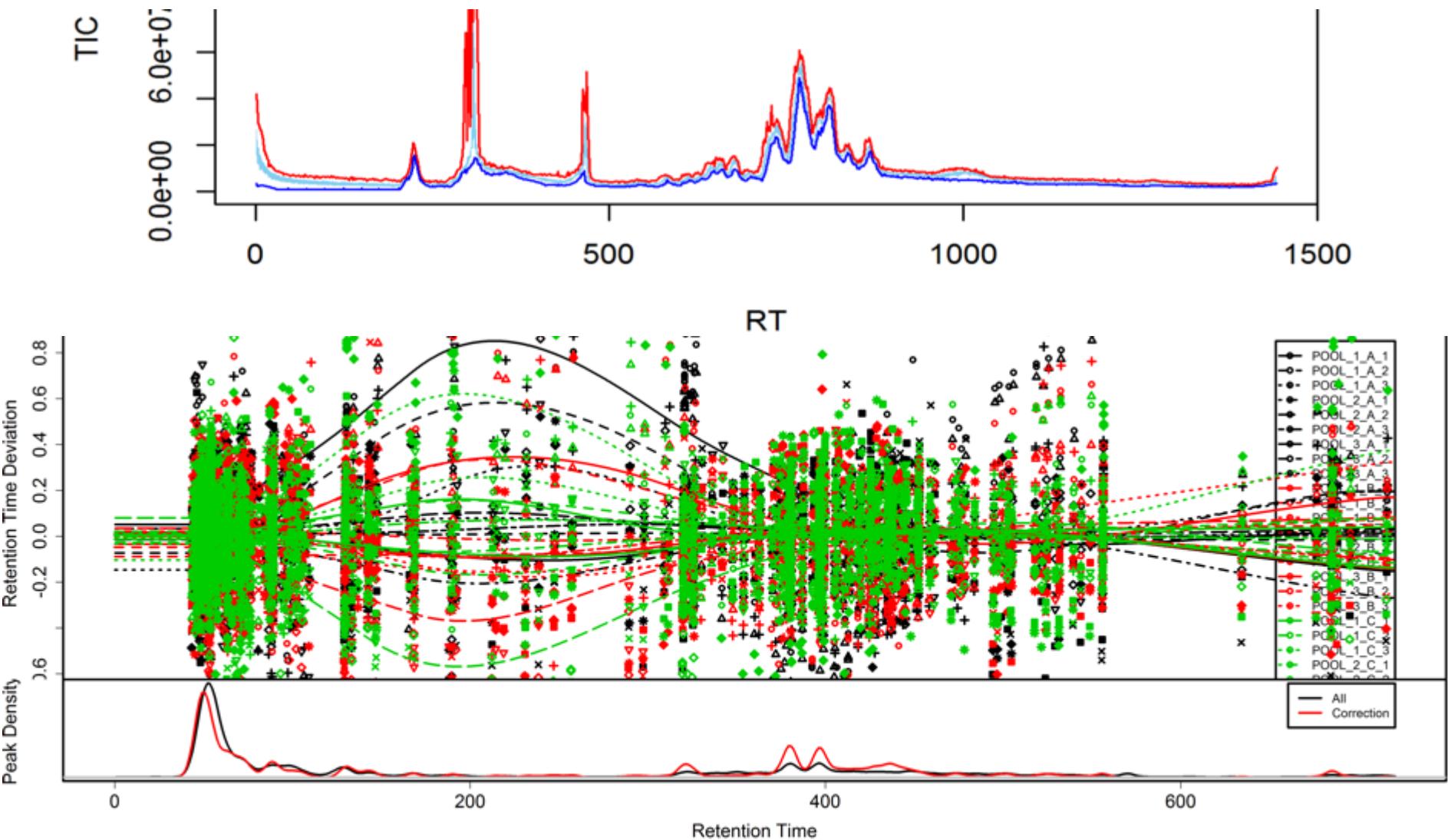
5. **Quantify**



Kenar et al., Mol. Cell. Prot., 2014, 13(1):348-59. doi: 10.1074/mcp.M113.031278

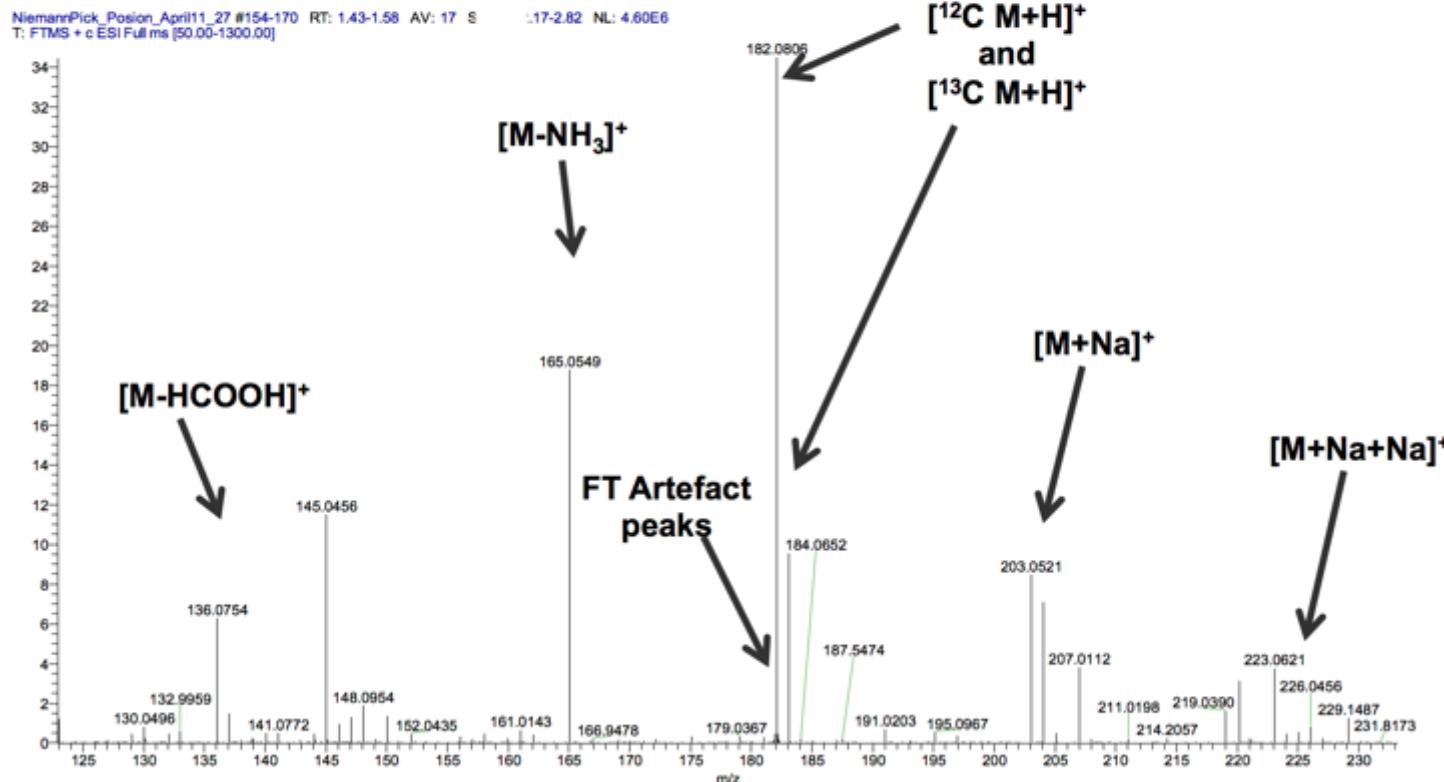
Oliver
Kohlbacher
University of
Tübingen

Grouping and retention time correction



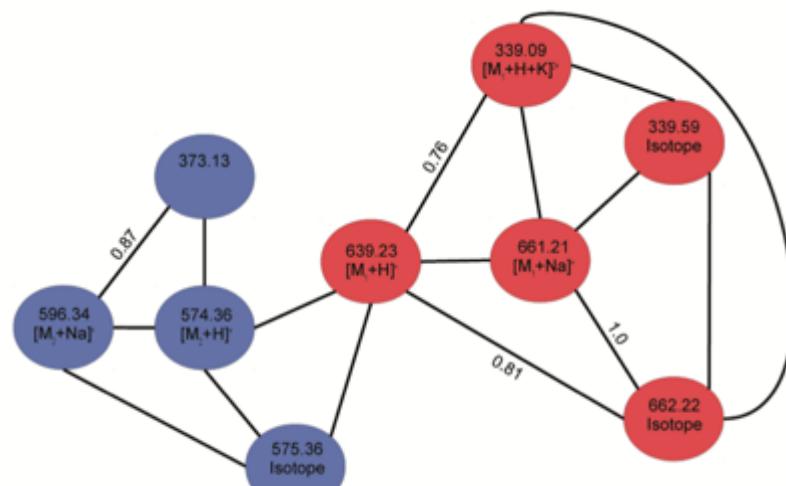
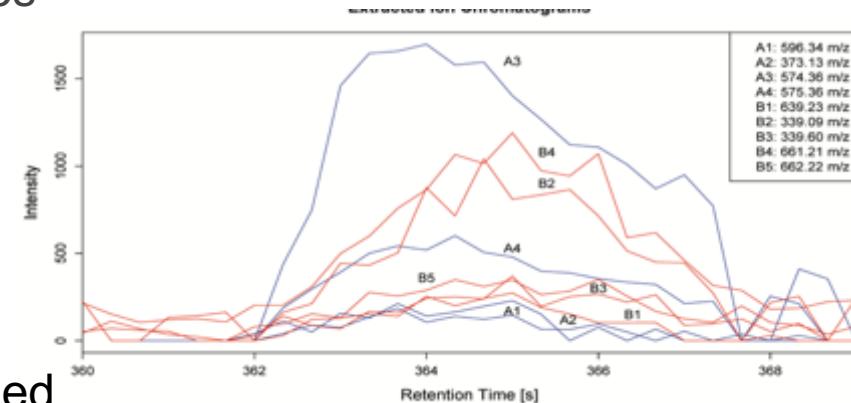
The complexity of data for Tyrosine

17 IONS OF DIFFERENT MASS AND SAME RETENTION TIME CREATED FROM A SINGLE METABOLITE – NEED TO DEFINE THE ION TYPE FOR IDENTIFICATION OR HIGH PROBABILITY OF FALSE POSITIVE IDENTIFICATION



Camera algorithms

1. **groupFWHM:** Using the most intense features the data is divided into rough retention time groups.
2. **Findisotops:** Look for C12/C13 isotope differences of m/z. Checks if intensity profile matches for $[M+]$ + to $[M+1]+$
3. **groupCorr:** Using the extracted ion chromatograms (EIC) for each feature are used to calculate *two types of correlation* between features:
 - Correlation across samples: CAS
 - Correlation within samples : CPSi
4. **Network relationship map** to build:
 - Nodes: features (peak)
 - Weighted Edges: The score (above a threshold)

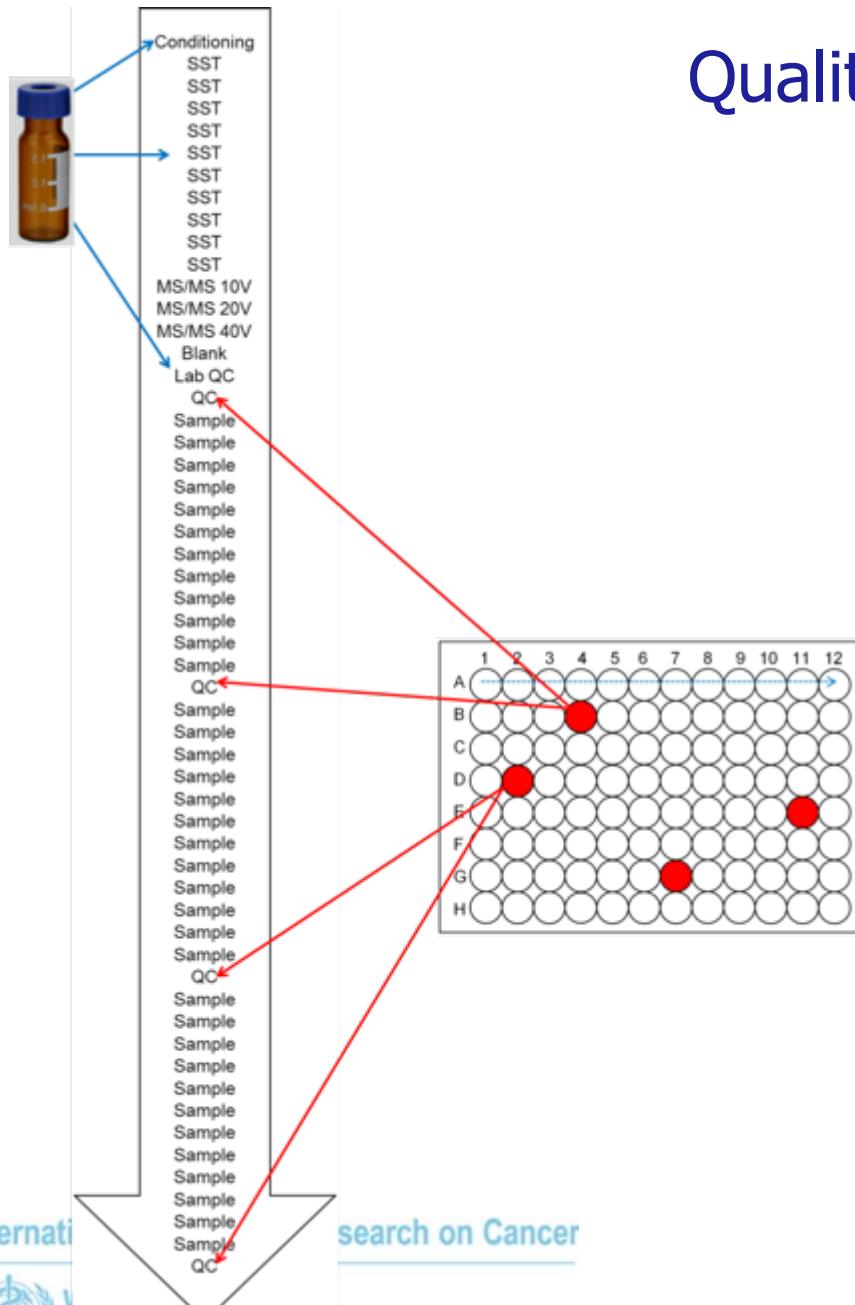


Peak tables and filtering, data clean-up

METABOLOMICS DATA

International Agency for Research on Cancer



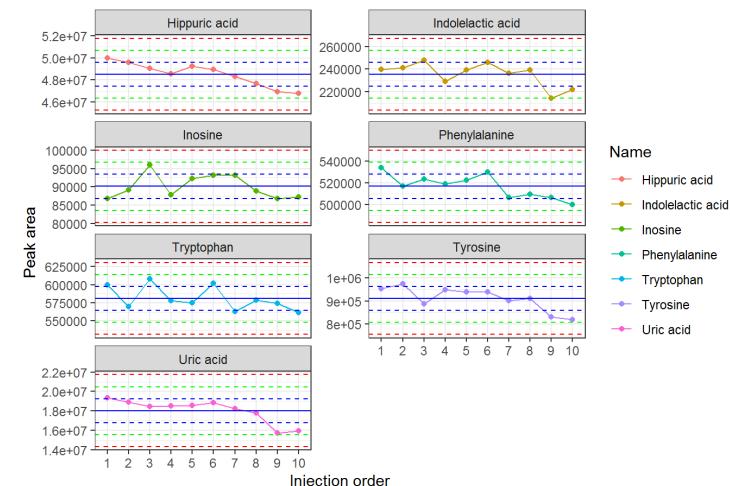


Quality control

QC on preselected metabolites

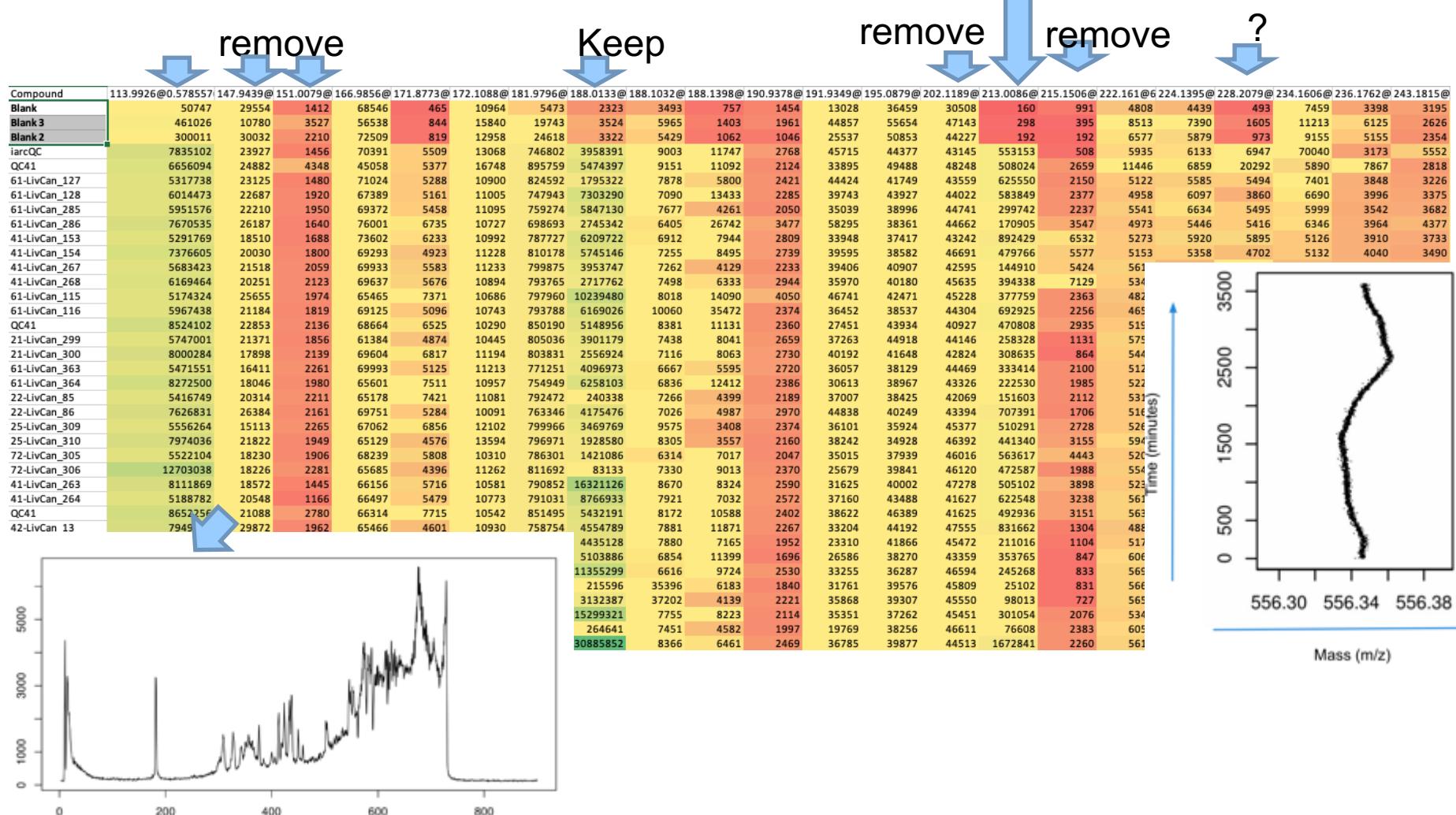
Metabolite	Mean (peak area)	Peak area RSD (%)
Cholesterol	2358180	6.3
Cortisol	128190	10.9
Creatine	3461812	9.0
Decanoylcarnitine	116241	9.8
Hippuric acid	626165	10.2
Hypoxanthine	1701776	8.7
Indole-3-acetic acid	1923894	8.4
Indolelactic acid	503305	7.7
Indolepropionic acid	399383	10.0
Inosine	36161	12.0
Kynurenone	1278834	11.7
Phenylalanine	27782412	8.9
Retinol	99445	10.3
Tryptophan	30302944	12.1
Tyrosine	2465081	14.4
Valine	43778692	3.3

Westgard rules: target metabolites in QC samples



Removing “signals” blanks

The primary purpose of blanks is to trace sources of artificially introduced contamination.



MaConDA

Method Blank: A blank prepared to represent the matrix as closely as possible. The method blank is prepared/extracted/digested and analysed exactly like the field samples.

Purpose: Assess contamination introduced during sample preparation

Instrument Blank: A blank analysed with field samples.

Purpose: Assess the presence or absence of instrument contamination.

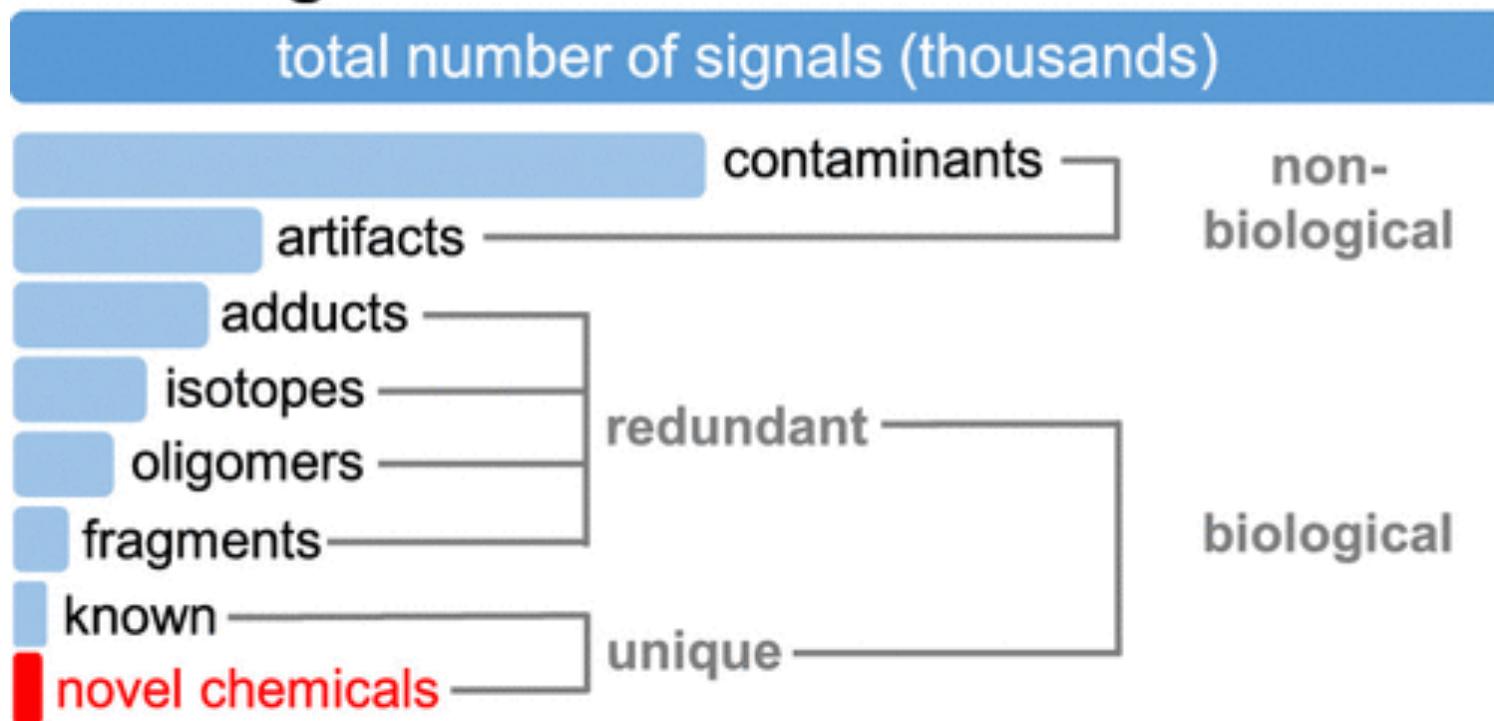
Frequency: Defined by the analytical method or at the analyst's discretion (e.g., after high concentration samples).

The screenshot shows the MaConDA website. At the top, there is a search bar with a dropdown menu set to "Name" and a "Search" button. Below the search bar is a navigation menu with links for Home, Search, Browse (which is highlighted), About, Downloads, and Contact. The main content area displays a table of mass spectrometry contaminants. The table has columns for ID, Exact mass, Name, Formula, and Type. The data rows are color-coded in shades of blue and white.

ID	Exact mass	Name	Formula	Type
CON00103	82.003075	Sodium acetate	C2H3O2Na	Solvent
CON00104	135.974808	Sodium trifluoroacetate	C2F3O2Na	Salt
CON00105	283.287506	Stearamide	C18H37NO	Slip agent
CON00106	284.271515	Stearic acid	C18H36O2	
CON00107	557.566650	Stearyl-palmityldimethylammonium chloride	C36H76NCl	Personal care products
CON00108	97.967384	Sulphuric acid	H2SO4	
CON00109	242.284775	Tetrabutylammonium (TBA)	C16H36N	Buffer
CON00110	518.131531	Tetradecamethylcycloheptasiloxane	C14H42Si7O7	Airborne Contaminant
CON00111	336.111115	Tributyl Tin Formate	C13H28O2Sn	Catalyst
CON00112	185.214355	Tributylamine	C12H27N	Solvent
CON00113	266.164703	Tributylphosphate	C12H27O4P	Plasticiser
CON00114	101.120445	Triethylamine	C6H15N	Buffer

Metabolomics Data reduction

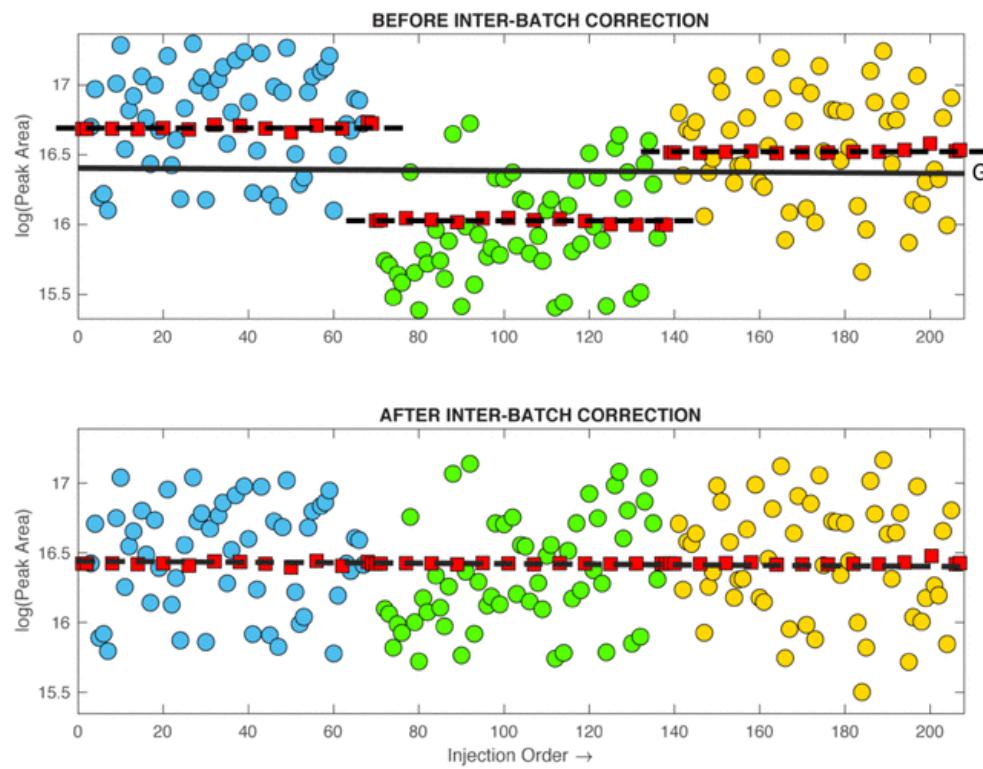
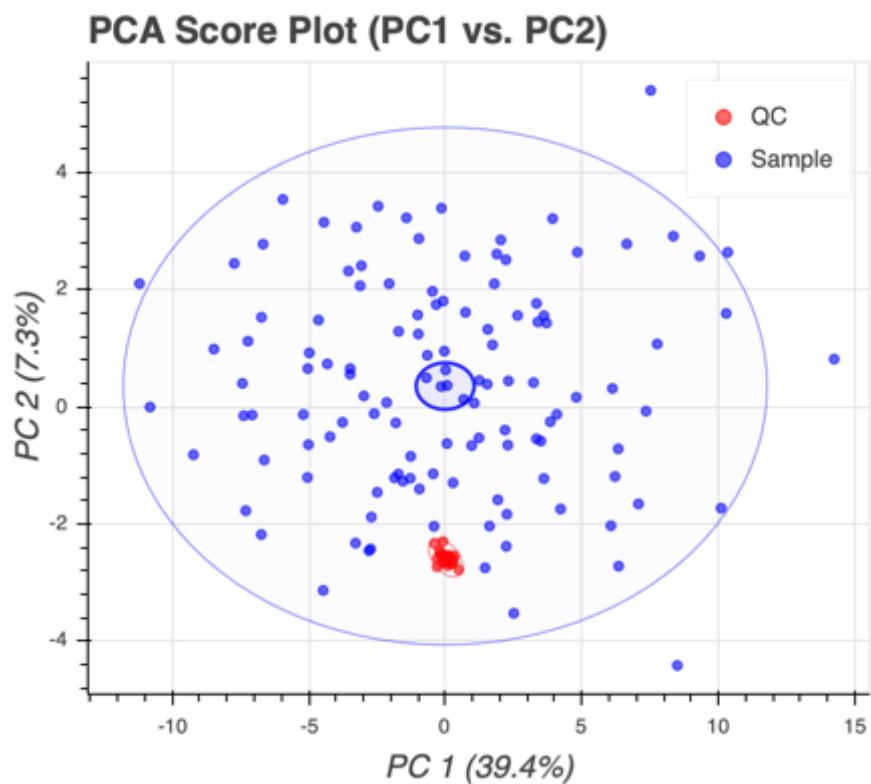
untargeted metabolomics - LC/MS data



<https://doi.org/10.1021/jacs.9b13198>

<https://yufree.github.io/pmd/>

QC and Batch correction



<https://doi.org/10.1007/s11306-018-1367-3>

Metabolite identification **METABOLOMICS DATA**

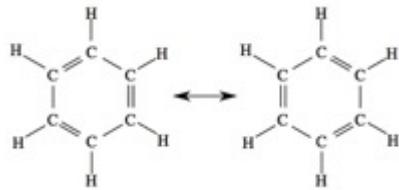
International Agency for Research on Cancer



Metabolome - Dark Matter

Dynamic Range and Instrument Sensitivity

Heterogeneity



multiple methods for sample prep and analysis

Chemical Lability

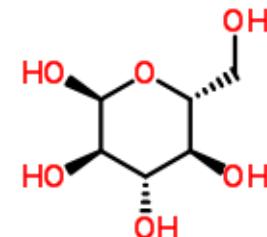
Keep samples cold,
minimum prep

Chemical Instability

Ionization Techniques and choice of MS



Isomers



Glucose (1 formula, 1 accurate mass)
385 structures

Adducts, fragments

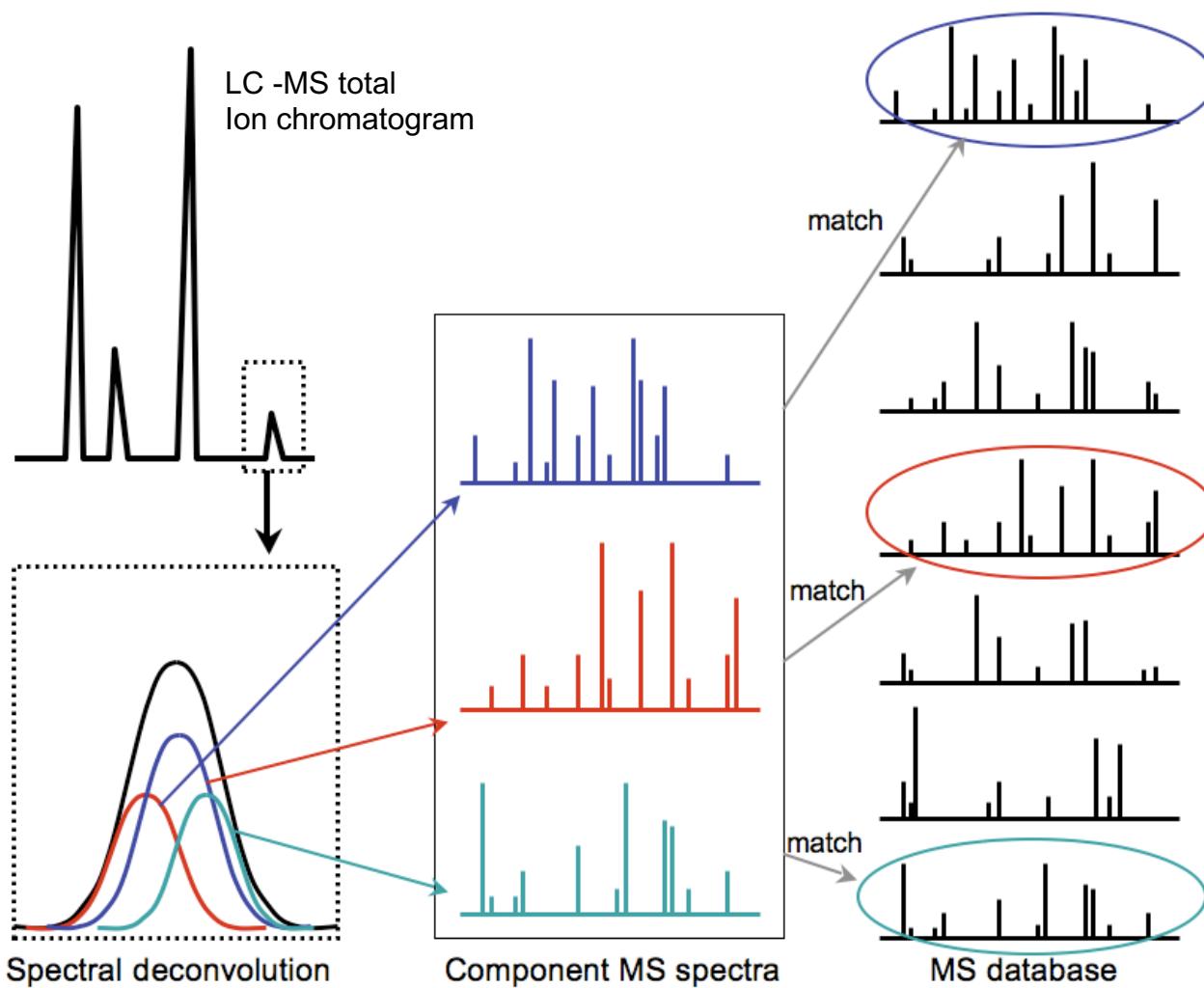
Mass spectrometers are machines for performing chemistry – Graham Cooks!

Database search

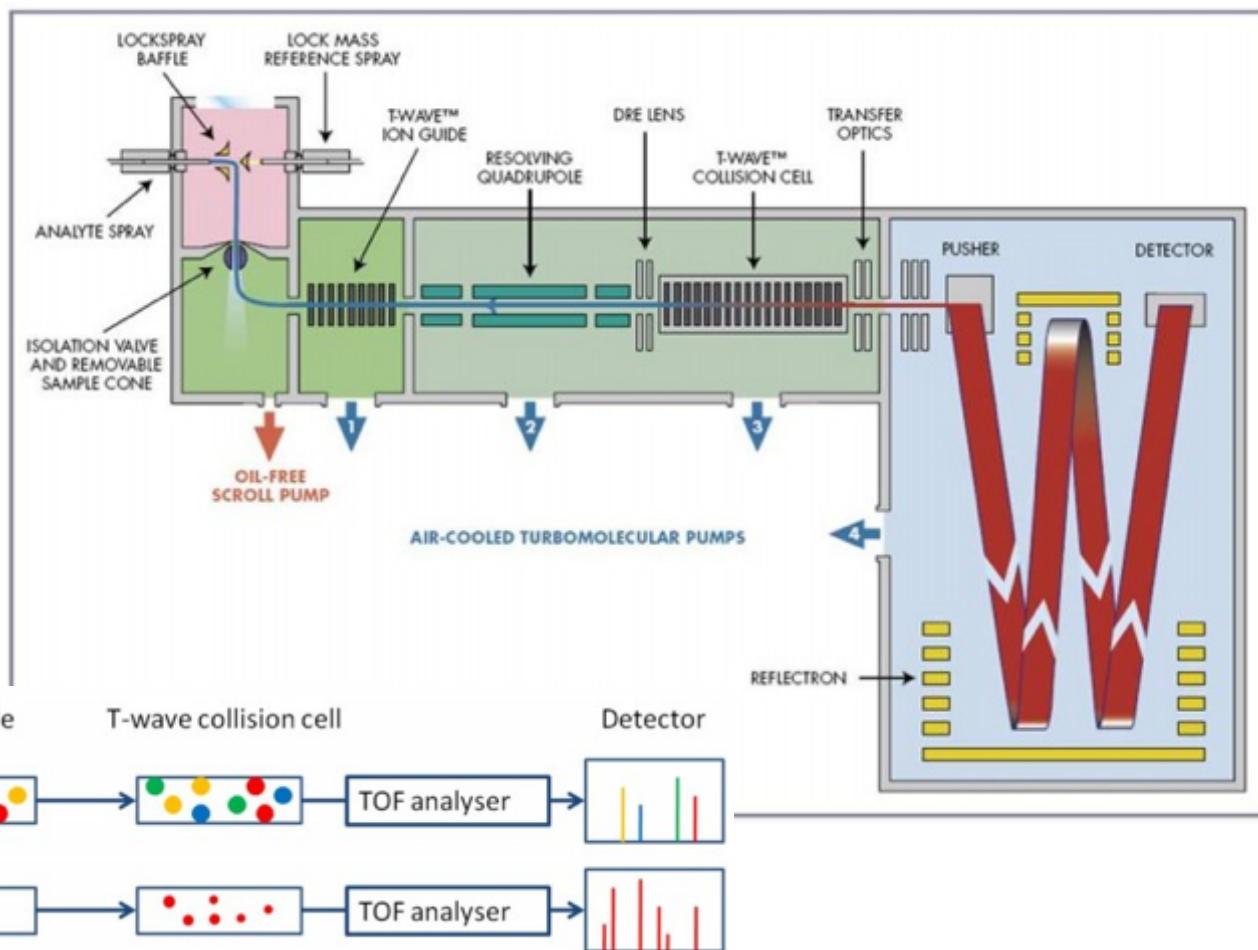
<https://hmdb.ca/>



Metabolite ID by LC-MS



Metabolite ID using MS/MS



Data tables and statical analysis

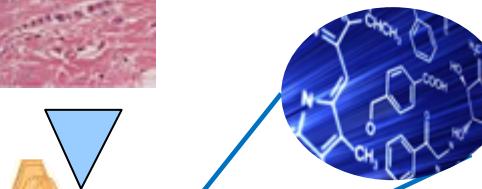
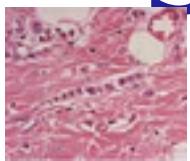
METABOLOMICS ANALYSIS

International Agency for Research on Cancer



World Health
Organization

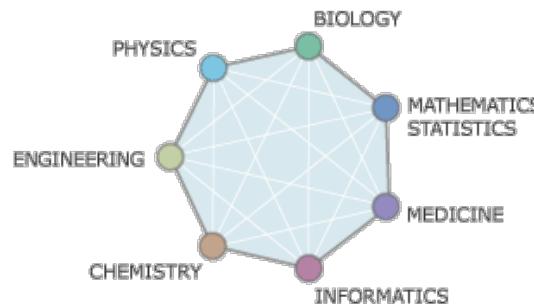
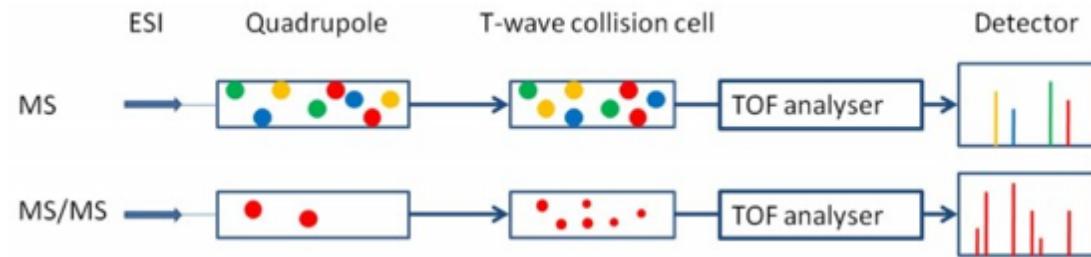
Snapshot



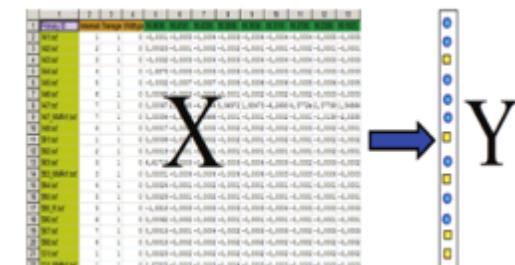
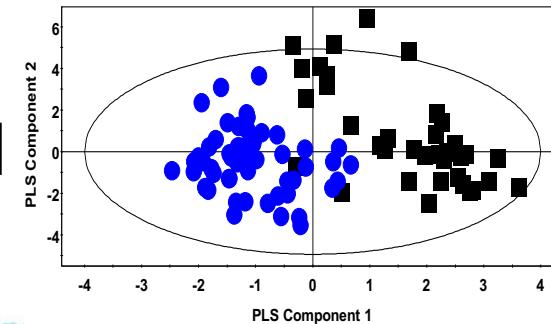
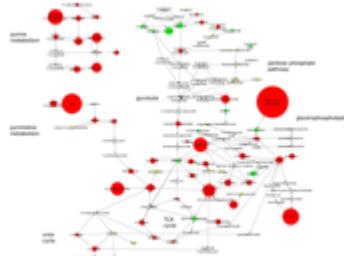
DNA &
Protein
Pellet

Aqueous
Layer

Organic
Layer



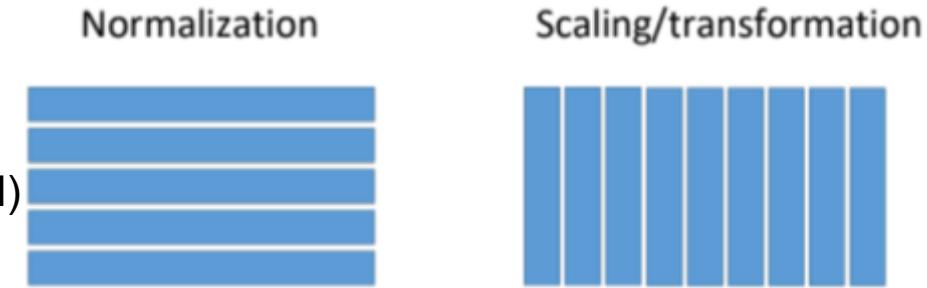
Oltvai-Barabasi, Science, Oct 02



International Agency for Research on Cancer

Preparing the data for statistical analysis

- **Data Normalization:** Making **each observation comparable** to others by **removing effects due to systematic biases**, which could originate from e.g., overall concentration differences or other systematic effects.
 - Total area normalization (TAN)
 - Total sum normalization (TSN)
 - Normalization to standard
 - Normalization to the major peak
 - Probabilistic Quotient Normalization (PQN)
 - Quantile normalization (QN)
- **Data Centring**
- **Data Scaling**
 - Scaling based on data dispersion
 - Autoscaling, Pareto scaling, range scaling, vast scaling
 - Scaling based on average value
 - Level scaling
- **Data Transformations**

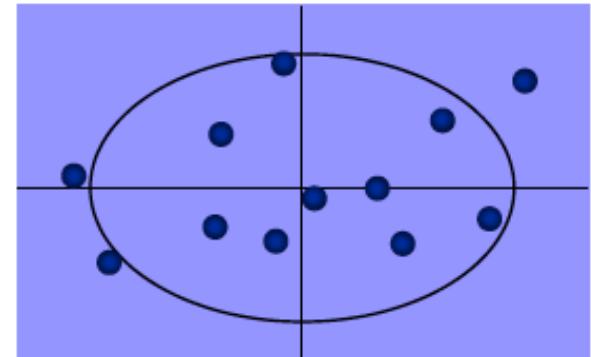
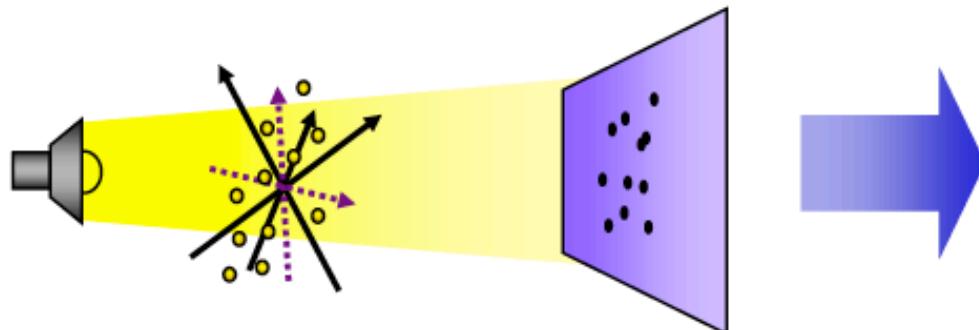


PCA and MVDA



Multivariate analysis by Projection

- Looks at ALL the variables together
- Avoids loss of information
- Finds underlying trends = “latent variables”
- More stable models



Metabolomics experimental database

EMBL-EBI About us Training Research Services

EMBL-EBI Hinxton

Homo sapiens

PUBLICATIONS

Individual variability in human blood metabolites identifies age-related differences.

Chaleckis R, Murakami I, Takada J, Kondoh H, Yanagida M.



Descriptors Protocols Samples Assays Metabolites Files

FTP Download

Aspera Download

Help

ISA METADATA

Download

search meta data

- a_mtbls265_NEG_mass_spectrometry.txt March 20 2017 14:39:37
- a_mtbls265_POS_mass_spectrometry.txt March 20 2017 14:39:37
- i_Investigation.txt December 04 2019 16:18:03
- m_mtbls265_NEG_mass_spectrometry_v2_maf.tsv March 20 2017 14:39:37
- m_mtbls265_POS_mass_spectrometry_v2_maf.tsv March 20 2017 14:39:37
- s_MTBL265.txt March 20 2017 14:39:37



RAW / DERIVED FILES

search raw files

- Person01_blood_youth_NEG.mzML February 24 2016 16:11:23



2012

International

World Health Organization

<https://www.ebi.ac.uk/metabolights/>



isatab

Risa

Metabolomics tools; PhenoMeNal

PhenoMeNal Gateway
Cloud Research Environment Portal

Home CRE App Library Help Login

App Library - Service Catalogue

App Library showcases our service catalogue listing 36 applications that are available via Galaxy workflows and Jupyter libraries through the Cloud Research Environment.

Functionality

- Preprocessing
- Annotation
- Post-processing
- Statistical Analysis
- Workflows
- Other Tools

Approaches

- Metabolomics
- Isotopic Labelling Analysis
- Lipidomics
- Glycomics

Instrument Data Types

- MS
- LC-MS
- GC-MS
- DI-MS
- CE-MS
- NMR
- IR
- Raman
- UV/VIS
- DAD

Search for Apps

Grid List

BATMAN
Bayesian Automated Metabolite Analyser for NMR spectra (BATMAN).



IPO
A Tool for automated Optimization of XCMS Parameters



Iso2Flux
Open source software for steady state ¹³C flux analysis



IsoDyn
"C++" program simulating the dynamics of metabolites and their isotopic isomers in central metabolic network using kinetic model



W4M LCMS matching
Annotation of MS peaks using matching on a spectra database.



metabomatching
metabomatching identifies metabolites using genetic spiking.



MetFrag
Command Line Interface for MassEins



MIDcor
"R"-program that corrects ¹³C mass spectrometric amounts of metabolites



W4M Metabolights Downloader
Metabolights downloader



metabolites
IMPACT FACTOR 3.303

MS data handling & processing
NMR data handling & processing
UV data
Ion species grouping & annotation
Statistics
Network analysis & biochemical pathways
R
Multifunctional workflows

The metaRbolomics Toolbox in Bioconductor and beyond

Volume 9 • Issue 10 | October 2019

MDPI mdpi.com/journal/metabolites ISSN 2218-1989

<https://portal.phenomenal-h2020.eu/>

International Agency for Research on Cancer



Acknowledgments

Biomarkers Group

Augustin Scalbert
Roland WEDEKIND
Vanessa Neveu
Genevieve Nicolas
Laure DOSSUS
Manon CAIRAT
Pekka Keski-Rahkonen
Nivonirina Robinot
David Achaintre
Sabina Rinaldi
Audrey Gicquiau
Mathilde His
Anne-Sophie NAVIONIS
Reza Salek
Adam Amara
Parisa Shahnazari
Pedro Ruas
Vanessa Neveu

Collaboration

Steffen Neumann
Fabien Jourdan
Francisco Couto
Justin van der Hooft
Kati Hanhineva
David Wishart



French National Institute for Health and Medical Research

French National Institute for Cancer Research

EPIC PIs

