

# Introduction to Metabolomics Platforms and Data Analysis

Dr. Katerina Kechris  
Professor  
Department of Biostatistics and Informatics  
University of Colorado Anschutz Medical Campus

# Outline

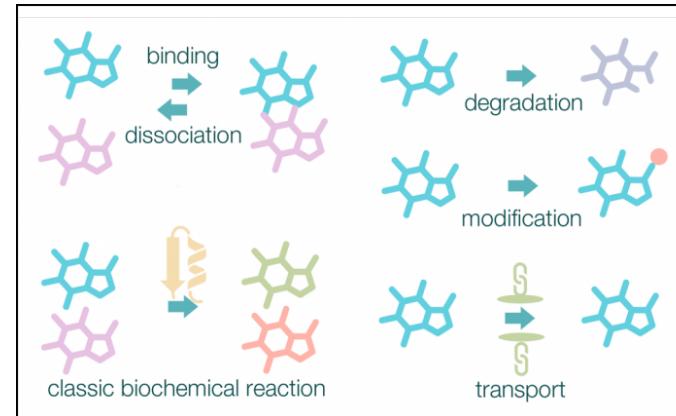
- Introduction
- Technology Types
- Mass Spectrometry
- Analysis Methods

# Introduction

- Metabolomics is the study of small molecules in biological or ecological samples
- Molecules can be identified utilizing analytical strategies
- Methods can be targeted or un-targeted
- Primary goal is to accurately identify and quantify levels of small molecules in samples

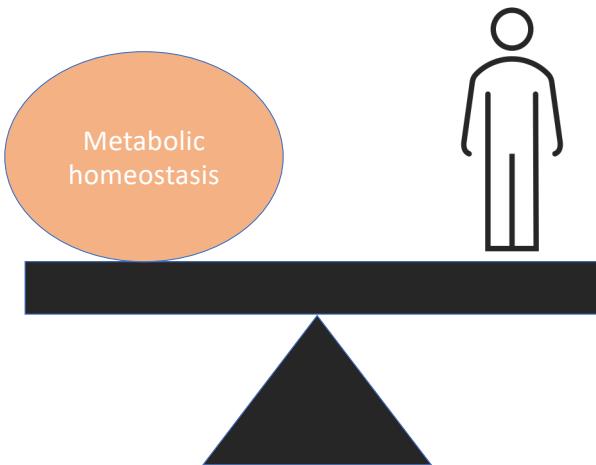
# Metabolism

- Series of chemical reactions
  - Enzymes
  - Catabolism/Anabolism

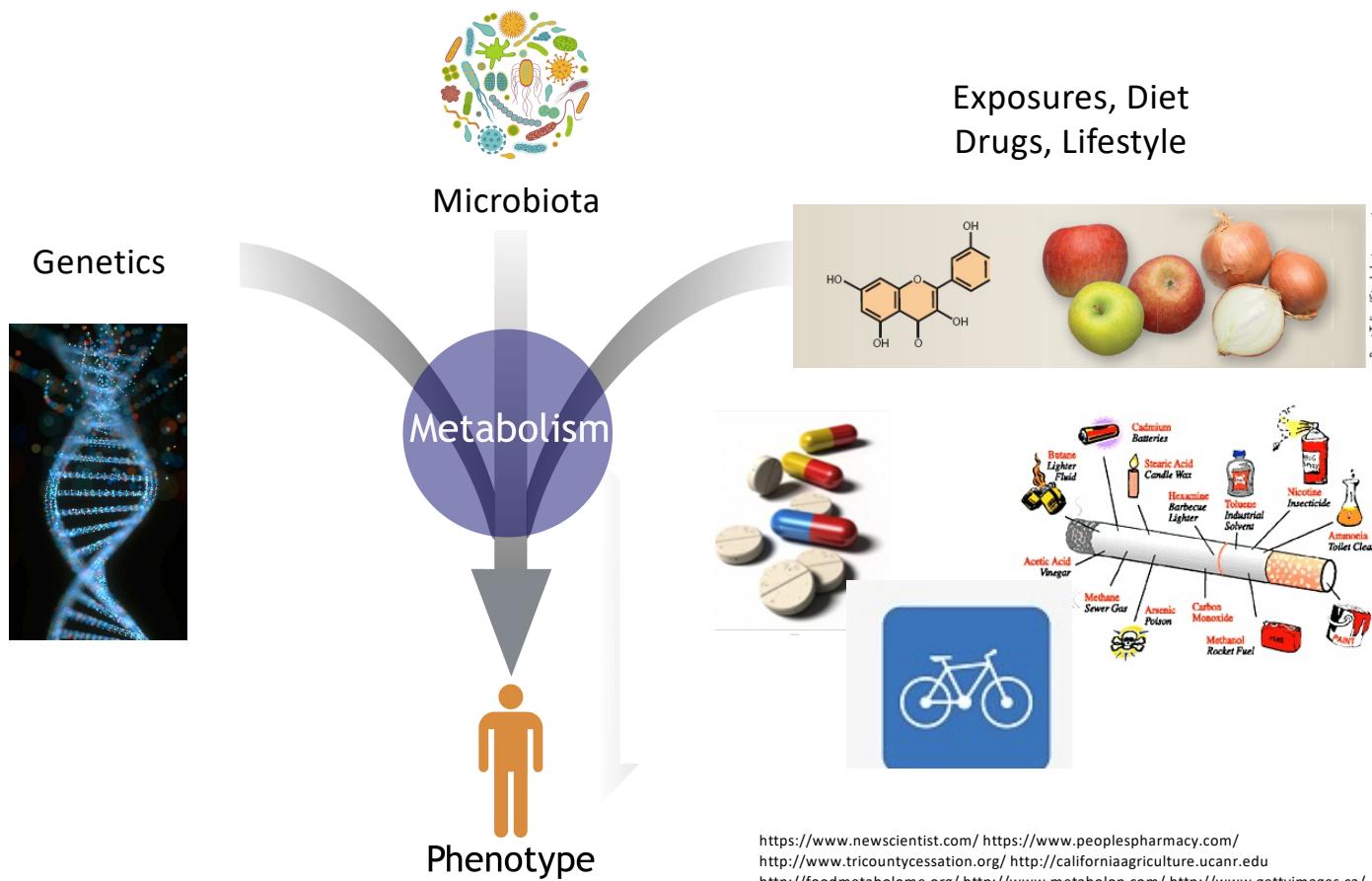


<https://www.ebi.ac.uk/>

- Homeostasis
  - Metabolites maintain balance

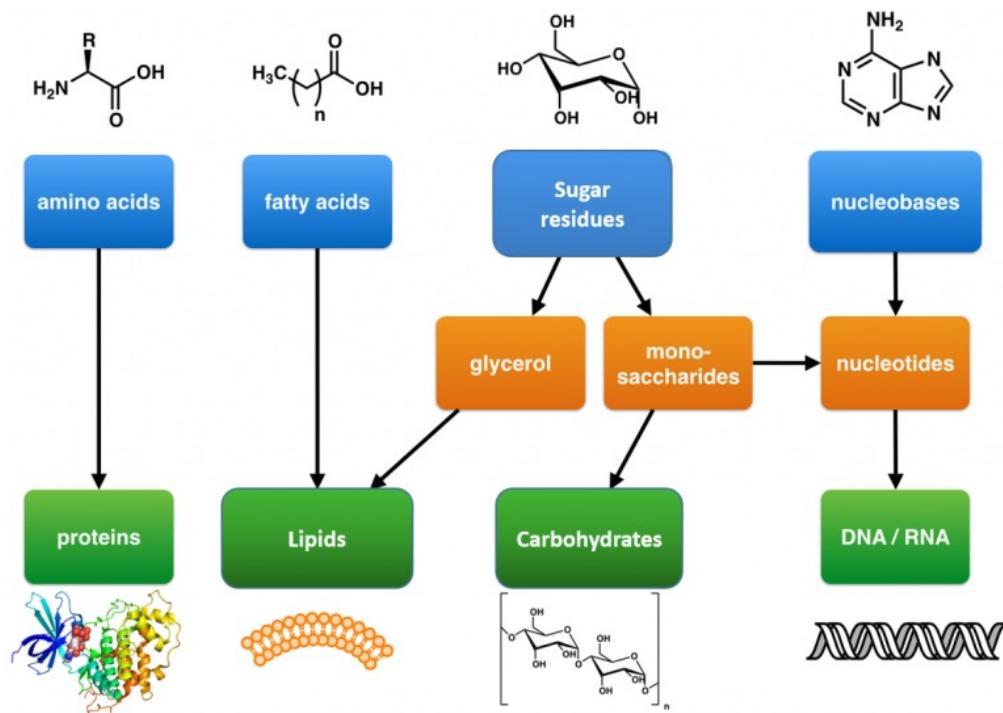


# Metabolism



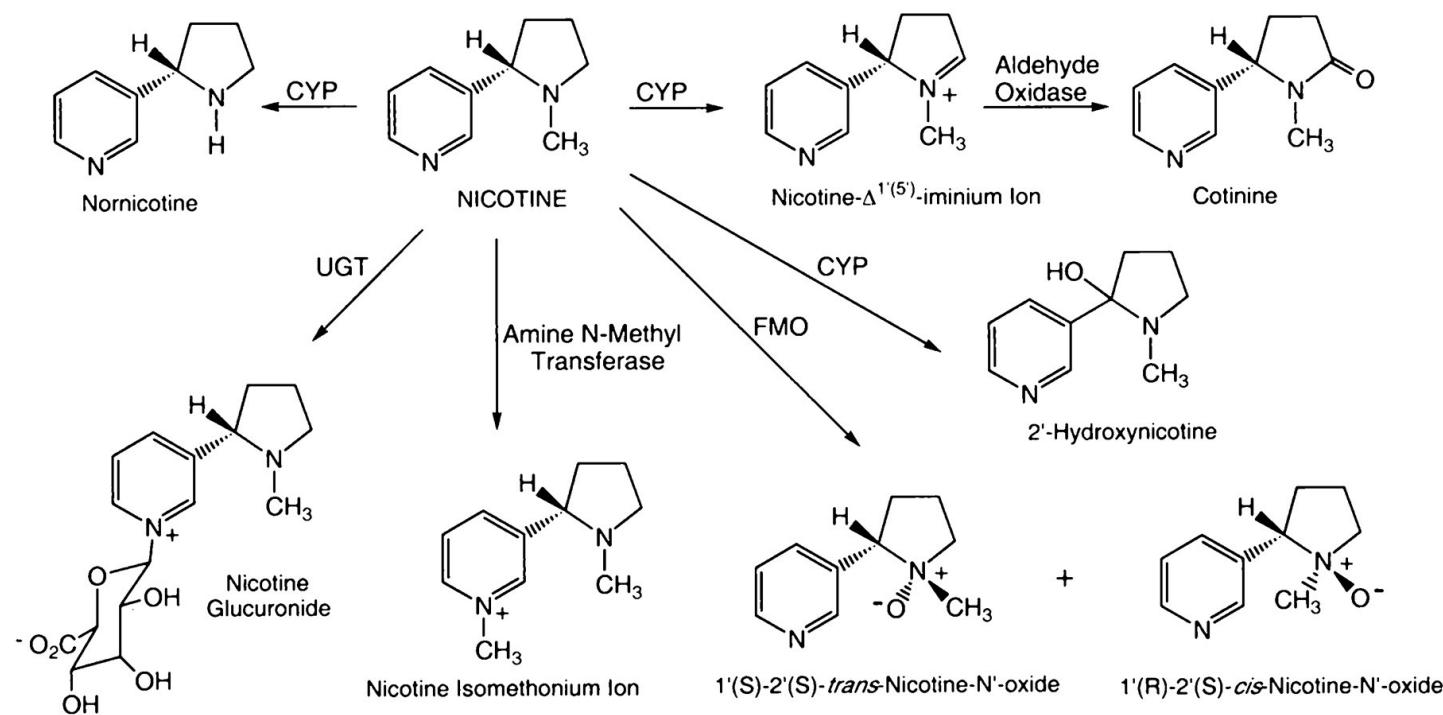
# Classes

## Small Molecules/Compounds/Metabolites



<https://wou.edu/chemistry/chapter-11-introduction-major-macromolecules/>

# Example: Nicotine Metabolism



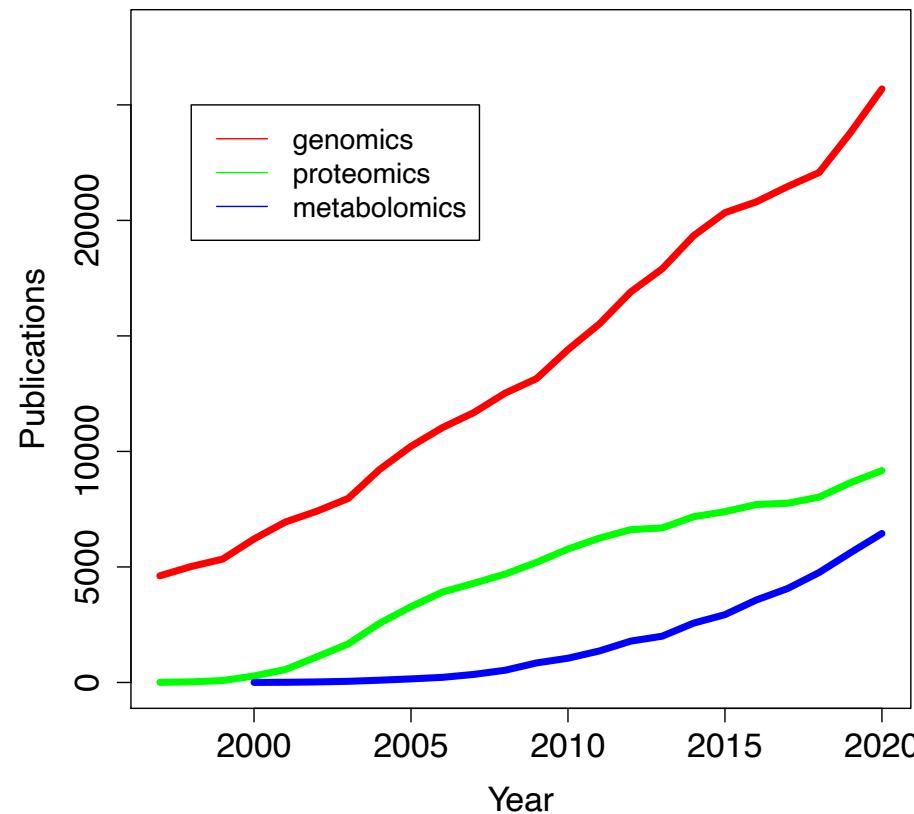
Nakajima & Yokoi (2005) Drug Metab Pharmacokinet

# Early History



- Linus Pauling and Noah Robinson
  - Profiled 50-200 metabolites early GC-MS
  - Profile disease, gender, and more (15K subjects)
  - 1000 ft stainless steel column
- 
- Term “metabolomics” coined
- 
- Wishart group (U of Alberta)
  - Human Metabolome Project
  - 2500 human, 1200 drug, and 3500 food metabolites

# ‘Omics Publications



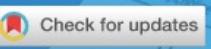
Source: PubMed 03/12/21, searched -omic\* in “TextWord”

# Translational Research

OPINION SPECIAL ISSUE: BIOMARKERS OF SUBSTANCE ABUSE | VOLUME 24, ISSUE 2,  
P197-205, FEBRUARY 01, 2018

Using Metabolomics to Investigate Biomarkers of Drug Addiction

Reza Ghanbari • Susan Sumner  

Published: February 01, 2018 • DOI: <https://doi.org/10.1016/j.molmed.2017.12.005> •  Check for updates

PDF [1 MB] Figures Save Share



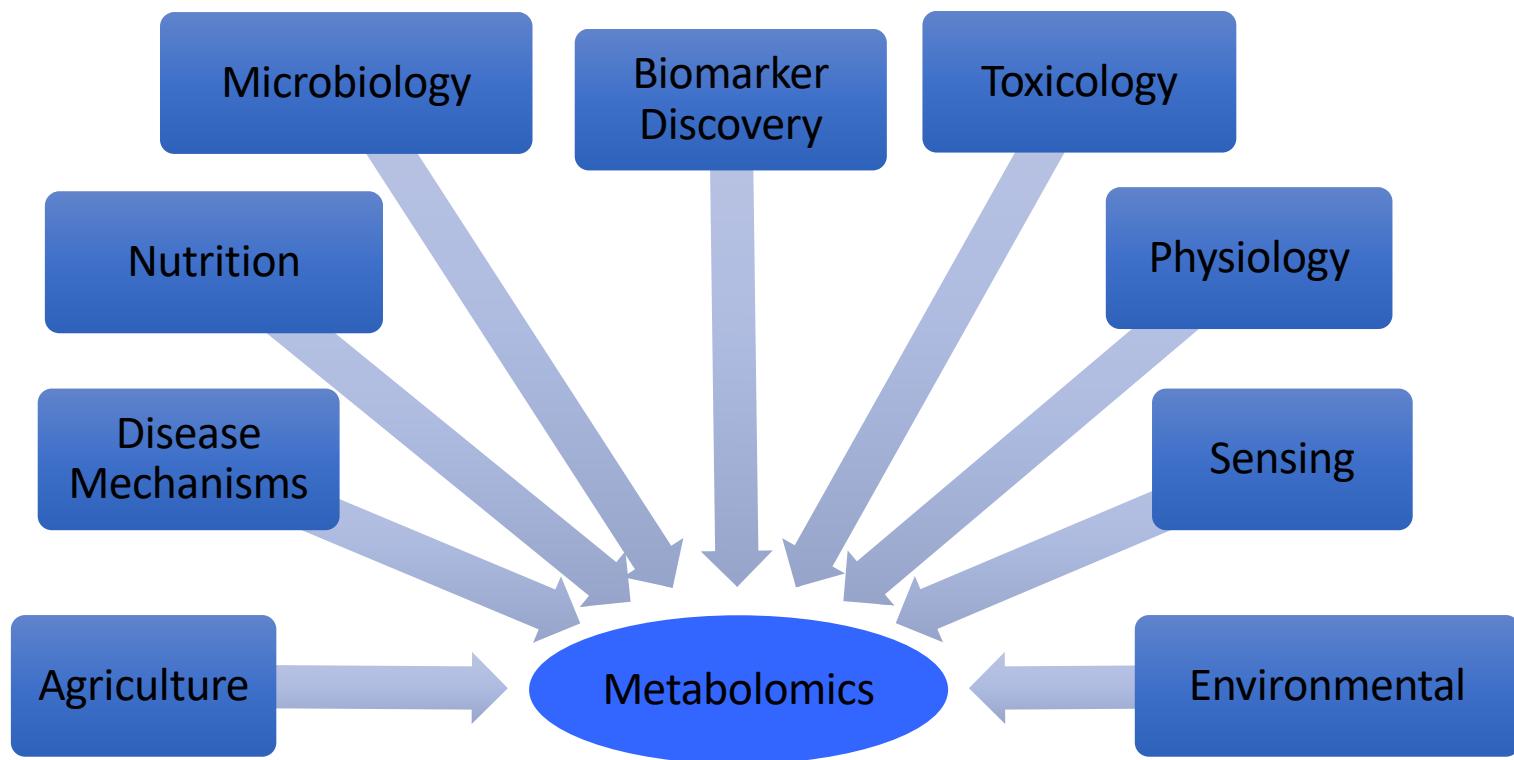
## Plasma metabolomic profiles enhance precision medicine for volunteers of normal health

Lining Guo<sup>a,1</sup>, Michael V. Milburn<sup>a</sup>, John A. Ryals<sup>a</sup>, Shaun C. Lonergan<sup>a</sup>, Matthew W. Mitchell<sup>a</sup>, Jacob E. Wulff<sup>a</sup>, Danny C. Alexander<sup>a</sup>, Anne M. Evans<sup>a</sup>, Brandi Bridgewater<sup>a</sup>, Luke Miller<sup>a</sup>, Manuel L. Gonzalez-Garay<sup>b</sup>, and C. Thomas Caskey<sup>c,1</sup>

<sup>a</sup>Metabolon, Inc., Durham, NC 27713; <sup>b</sup>Center for Molecular Imaging, Division of Genomics and Bioinformatics, The Brown Foundation Institute of Molecular Medicine, University of Texas Health Science Center, Houston, TX 77030; and <sup>c</sup>Department of Molecular and Human Genetics, Baylor College of Medicine, Houston, TX 77030

*Science, 2015*

# Key Applications



# Types of Questions

- **Identification:** What small molecules are present?
- **Association:** Can metabolites be used as biomarkers? Are metabolite levels associated with disease status/spectrum? physiological state? diet? environmental exposures?
- **Prediction:** Do small molecules have predictive power for diagnosis? prognosis? precision medicine?
- **Mechanisms:** What are the relevant metabolic mechanisms?
- **Genetics:** Is there variability due to genetic background?
- **Microbial effects:** What is the role of the microbiome on metabolism?
- **Metabolic flux:** What is the rate of turnover of molecules through a metabolic pathway?

[PLoS One.](#) 2017; 12(6): e0178281.

PMCID: PMC5456044

Published online 2017 Jun 2. doi: [10.1371/journal.pone.0178281](https://doi.org/10.1371/journal.pone.0178281)

PMID: [28575117](#)

## Gene and metabolite time-course response to cigarette smoking in mouse lung and plasma

Mikaela A. Miller,<sup>1</sup> Thomas Danhorn,<sup>2</sup> Charmion I. Cruickshank-Quinn,<sup>3</sup> Sonia M. Leach,<sup>2</sup> Sean Jacobson,<sup>4</sup> Matthew J. Strand,<sup>5</sup> Nichole A. Reisdorph,<sup>3</sup> Russell P. Bowler,<sup>4,\*</sup> Irina Petrache,<sup>4,6,\*</sup> and Katerina Kechris<sup>1,\*</sup>

# Experimental Design

	acute	subchronic			chronic	
	1 day	7 days	1 month	3 months	6 months	9 months
Air Control						
	x5	x5	x5	x5	x5	x5
Cigarette Smoke Exposed						
	x5	x5	x5	x5	x5	x5
Smoking Cessation						
						x5

# Altered Pathways

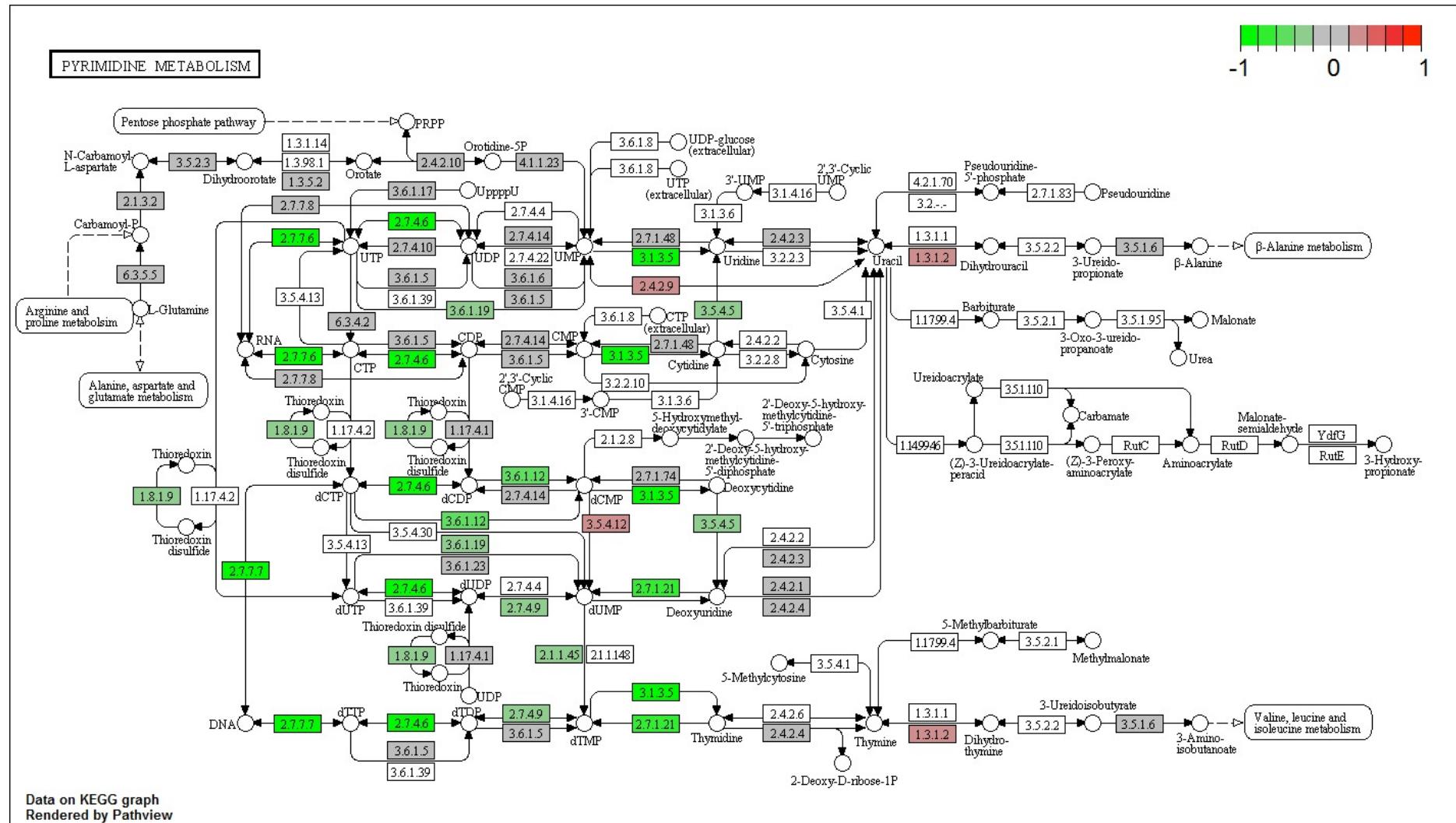
	1 day	7 days	3 months	6 months	9 months	Cessation
Oxidative phosphorylation				✓		
Parkinson's disease						✓
Pyrimidine metabolism					✓	✓
Leishmaniasis		✓	✓	✓		✓
Fc gamma R-mediated phagocytosis	✓	✓	✓	✓		✓
Inositol phosphate metabolism					✓	
Gap junction	✓	✓	✓			
Pathways in cancer	✓	✓	✓	✓		✓
Phosphatidylinositol (PI) signaling system	✓	✓	✓	✓	✓	✓

Checkmarks - perturbation in metabolomic data

Pink - enrichment for upregulated genes in CS-exposed mice (FDR≤0.10)

Blue - enrichment for downregulation in CS-exposed mice (FDR≤0.10)

Shaded fields under the “Cessation” column persisted following cessation when compared to AC mice



# Outline

- Introduction
- Technology Types
- Mass Spectrometry
- Analysis Methods

# Most Common Sample Types

- Plasma
  - contains fibrinogen (normal clotting of blood) & albumin (keeps fluid in bloodstream)
- Serum
  - plasma minus clotting factors and blood cells
- Urine
- But anything possible (saliva, stool, tissues, etc.)
- Considerations
  - fasting
  - time of day
  - hemolysis
  - freeze/thaw

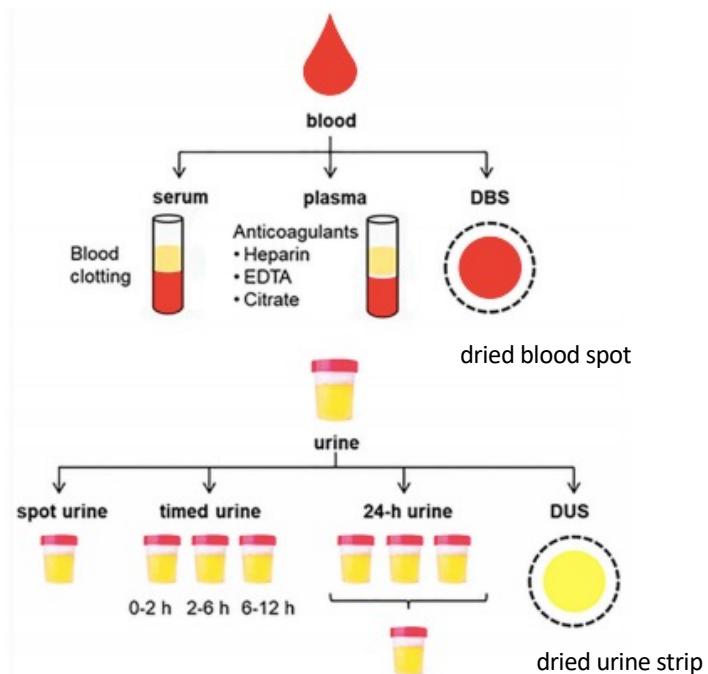
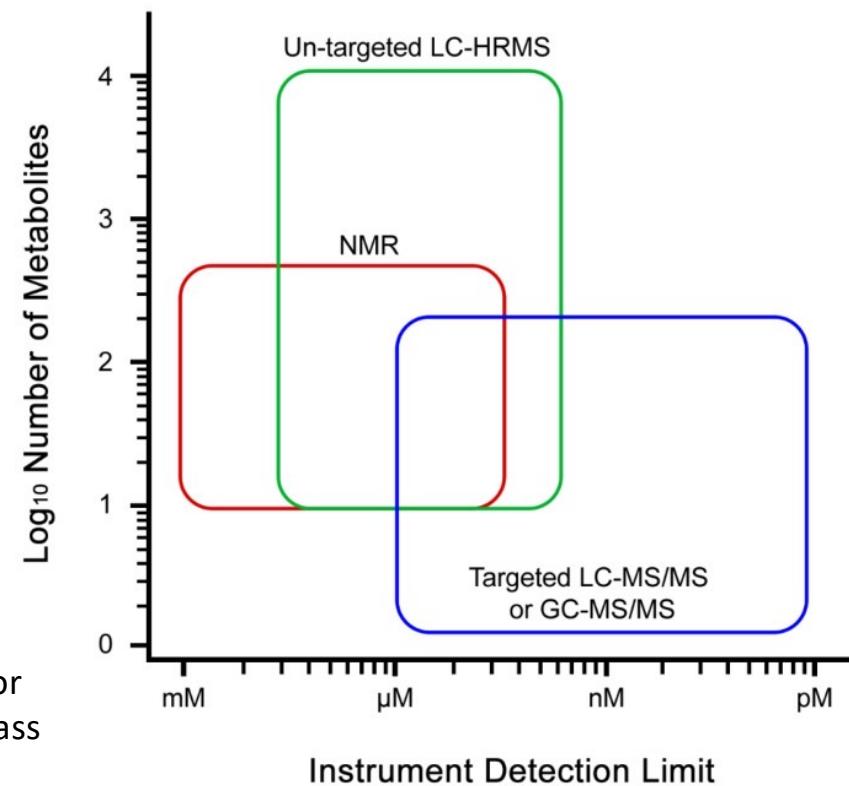


Figure 2. Summary of sampling alternatives for blood and urine collection.

González-Domínguez et al (2020) Metabolites

# Technologies

Liquid-Chromatography (LC) or Gas-Chromatography (GC) Mass Spectrometry (MS); Nuclear Magnetic Resonance (NMR)

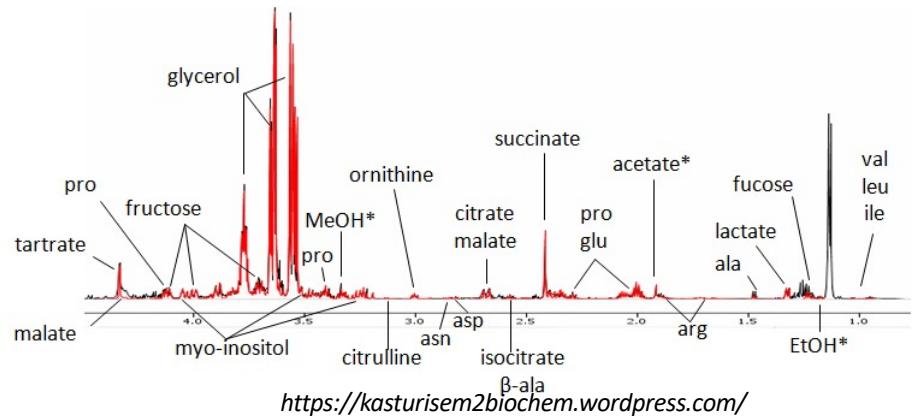


Pinu et al., (2019) *Metabolites*, [agilent.com](http://agilent.com), [bruker.com](http://bruker.com)



# NMR

- Sometimes referred to as “metabonomics”
- No separation needed & can handle intact tissues, non-destructive
- High analytical reproducibility, simple sample prep
- Absolute, quantitative data
- Also supports imaging & real-time metabolite profiling of living cells, metabolic flux analysis
- But less sensitive, and for a limited number of metabolites



# Mass Spectrometry Technologies

- Gas-Chromatography (GC)-MS
  - Uses fragmentation to identify molecules
- Liquid-Chromatography (LC)-MS
  - Can detect intact compounds and fragments but is destructive (sample is used)
  - Most common - ultra-high pressure (UHPLC)

# Types

## Targeted

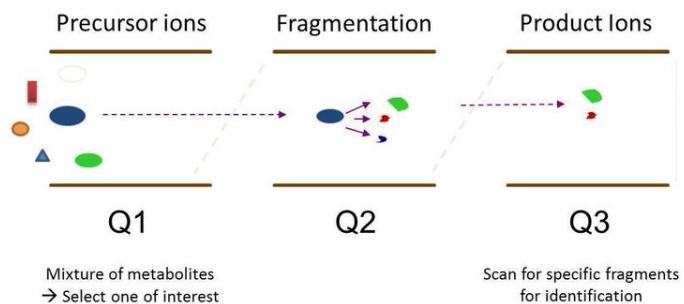
- Groups of characterized and annotated metabolites
- Hypothesis testing
- Known identities
- Absolute quantitation

## Untargeted

- Comprehensive coverage (limited only by techniques)
- Hypothesis generation
- Predicted identification
- Relative quantitation

# Targeted Analysis

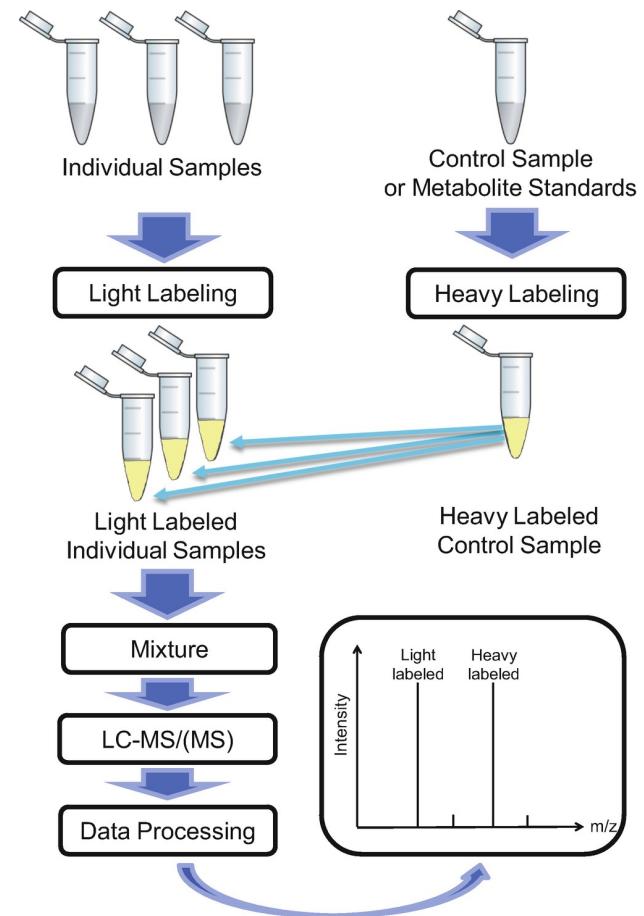
- Multiple reaction monitoring (MRM)
  - Parent mass of a metabolite of interest selected



- Use internal standards for identification and quantitation

# Internal Standards

- Light labeling = endogenous metabolites
- Use stable isotopes as heavy standards
  - No decay over time
  - No safety concerns (radioactivity)
- Spike in heavy standards at known concentration
- Measure relative intensities (peak areas) between light and heavy, then can back calculate concentration of endogenous metabolites



Zhao & Li (2021) In: Hu S. (eds) Cancer Metabolomics., vol 1280. Springer

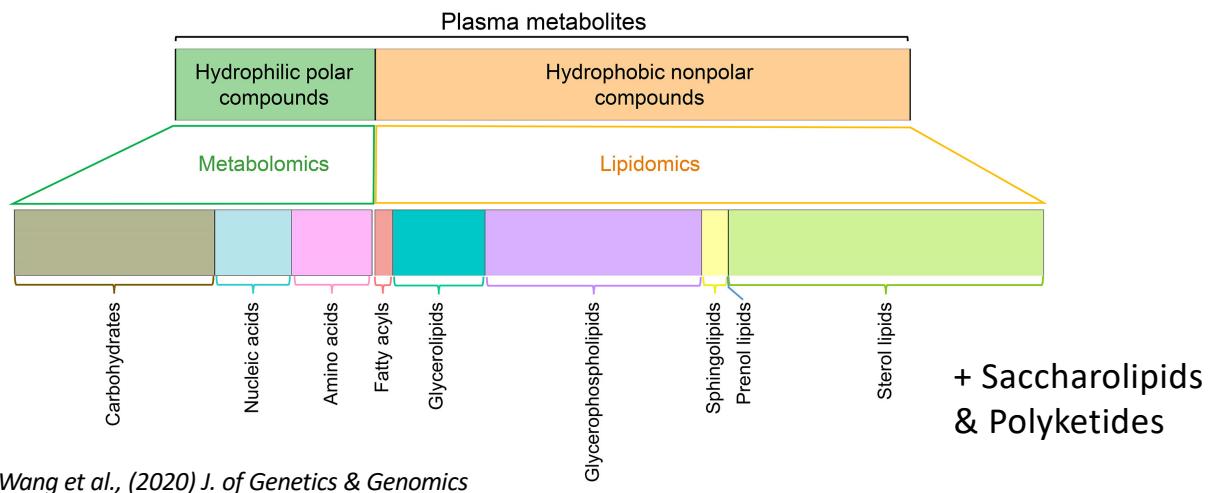
# Untargeted Analysis

- + Allows for novel discovery and comprehensive coverage
- Challenges include
  - time to process extensive amounts of raw data
  - difficulties in identifying unknown small molecules
  - coverage of the platform
  - bias towards detection of high-abundance molecules

Can be followed by targeted analysis to validate top candidates

# Lipidomics

- Profiling lipids requires own sample preparation, analytical protocol & data processing
- Within each lipid class, many variations of different length of carbon chains & number of double bonds
- Leads to particular patterns in mass spectra data, leveraged in data processing

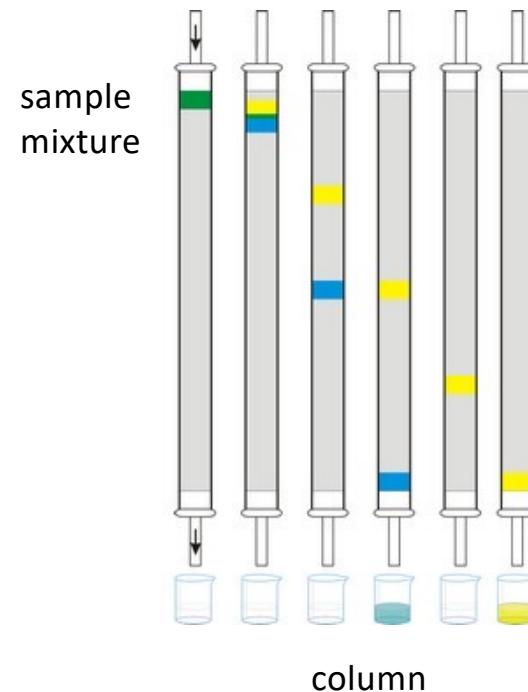


# Outline

- Introduction
- Technology Types
- Mass Spectrometry
- Analysis Methods

# Sample Preparation - Chromatography

- No digestion as in proteomics (to obtain peptides)
- Chromatography
  - Spreads out analytes over time
  - Less likely that ionization capacity overcome by large quantities of analyte
  - Two phases
    - Stationary Phase – analyte separation
    - Mobile Phase – carries analytes through chromatographic column to the mass spectrometer
  - Mobile phase: gas or liquid



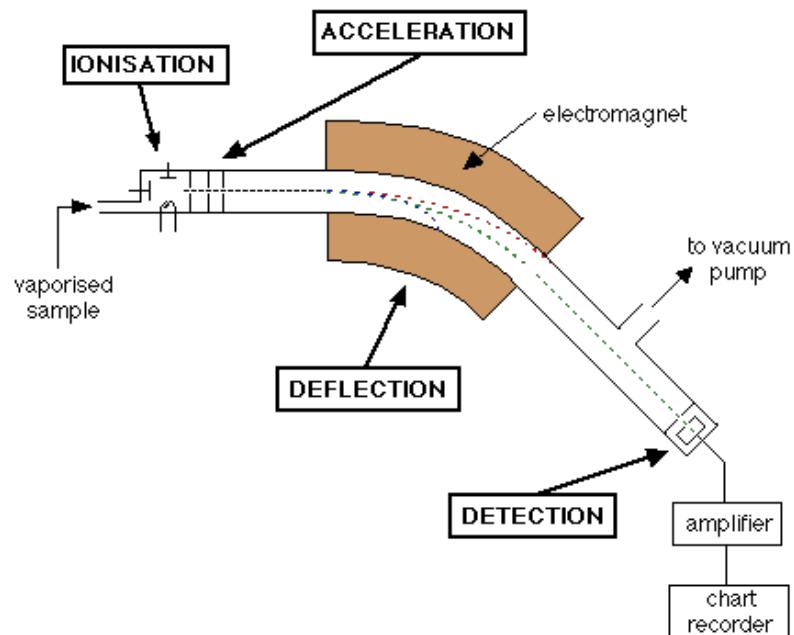
# Mass Spectrometry Basic Principles

Atoms and molecules can be deflected by magnetic fields if ionized (electrically charged)

1. Ionization – atom/molecule acquiring negative/positive charge  
*Targeted:* Only ions of preselected interest sent to the detector  
*Untargeted:* Reads all ions across a wide m/z range
2. Acceleration – ions are accelerated
3. Deflection – ions are deflected by magnetic field
4. Detection – beam of ions are detected electrically

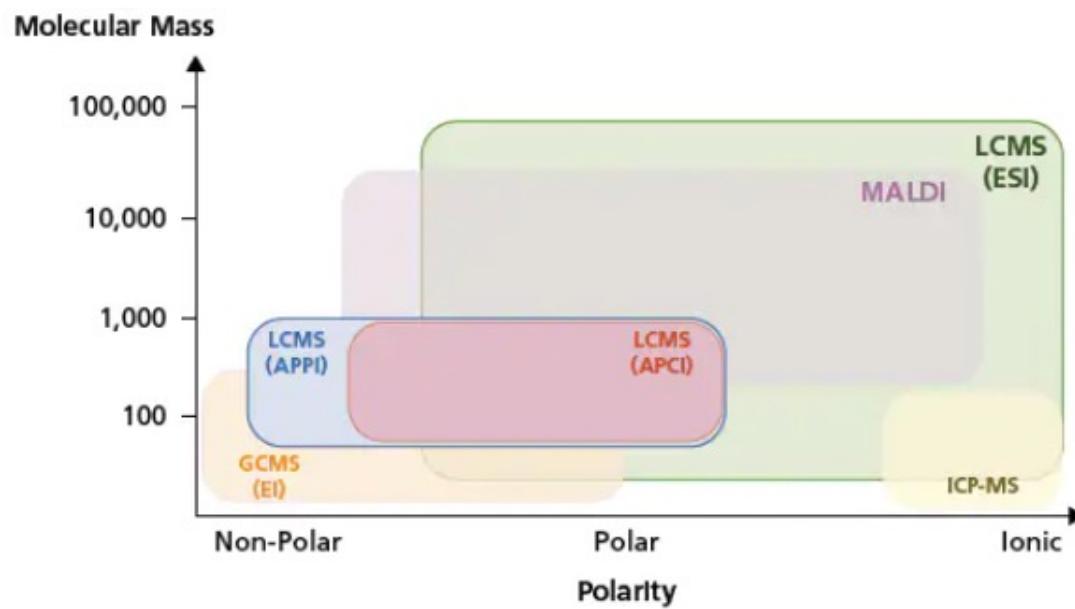
*Types :* Time of Flight (TOF), orbitrap

*High Resolution:* 4 decimal places offers molecular formula determination



# Ionization Step

During MS, analytes are ionized:

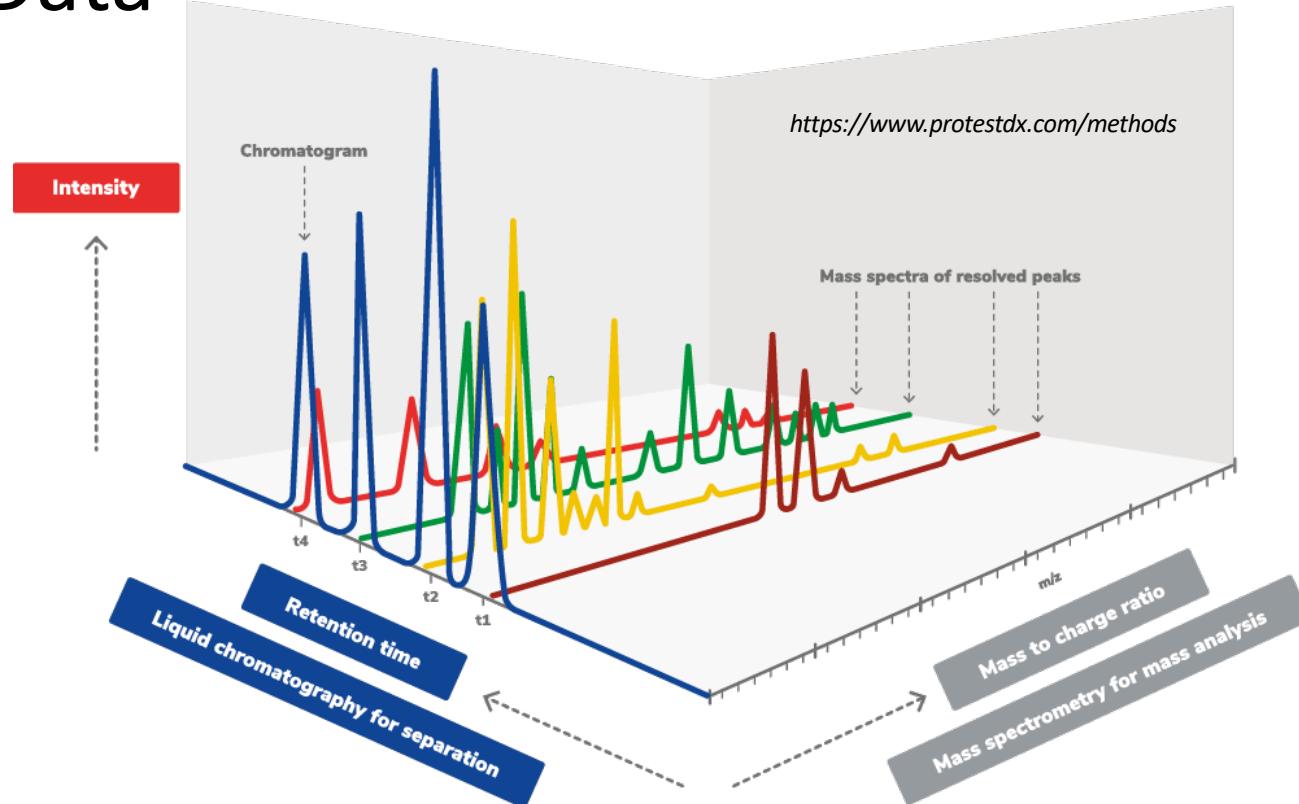


Electrospray  
ionization (ESI)  
(most popular)

Matrix-Assisted  
Laser Desorption  
Ionization (MALDI)

Atmospheric  
Pressure Chemical  
Ionization (APCI)

# Raw Data



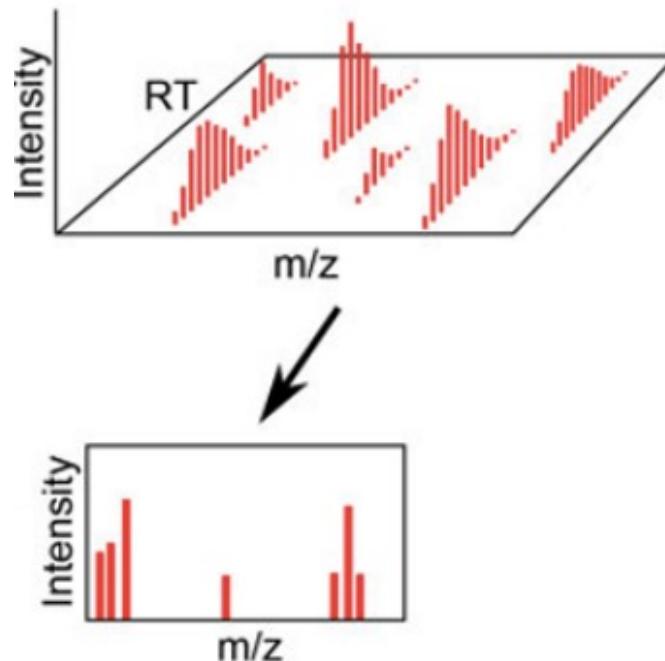
**Retention Time (RT):** time for solute to pass through chromatography column, time from injection to detection.

**MZ:** mass to charge ratio of compound ( $m/z$ )

# Raw Data

Compound identified based on RT and MZ

For fixed RT: m/z ratio (x-axis) by relative current (y-axis: relative abundance or intensity)



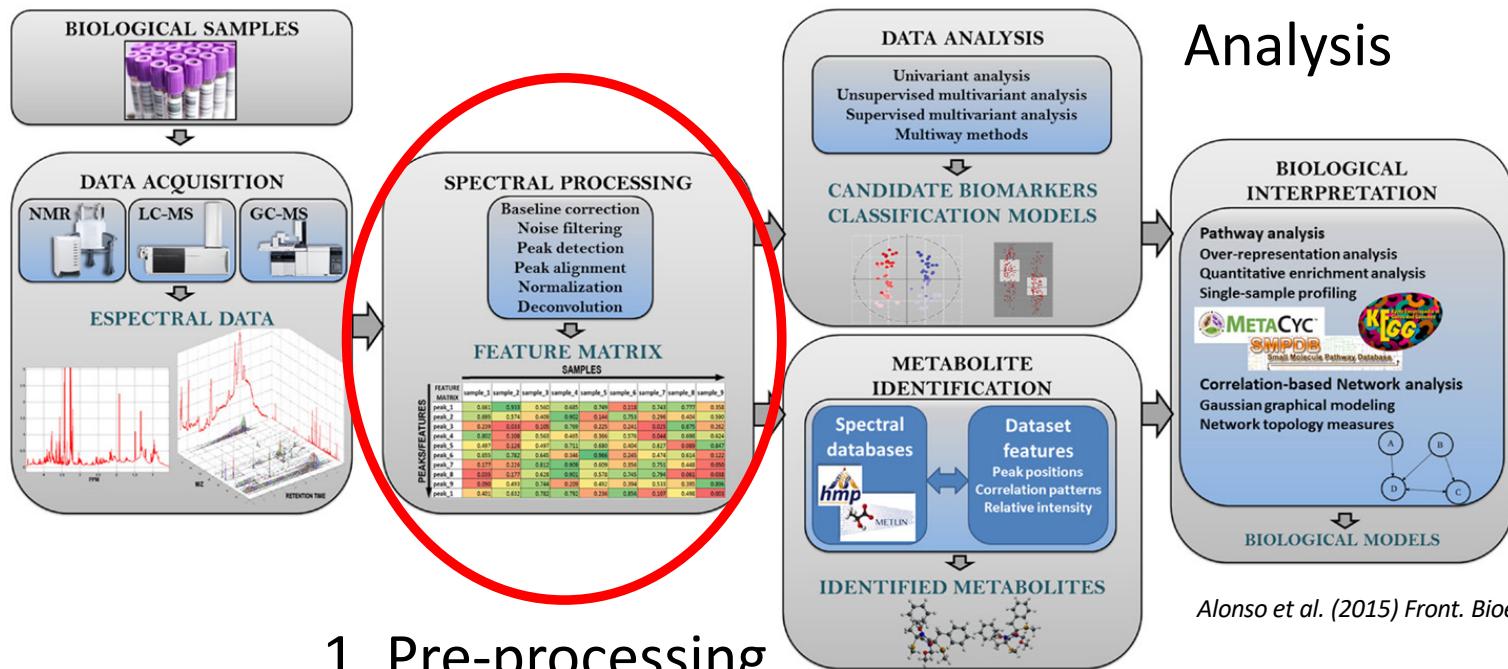
# Raw Data

- May be very large files (& binary)
- Often proprietary (only with commercial software)
- Open source formats
  - ANDI NetCDF, mzXML, mzData
  - Converters available (msConvert, FileConverter, etc)

# Outline

- Introduction
- Technology Types
- Mass Spectrometry
- Analysis Methods

# Pipeline

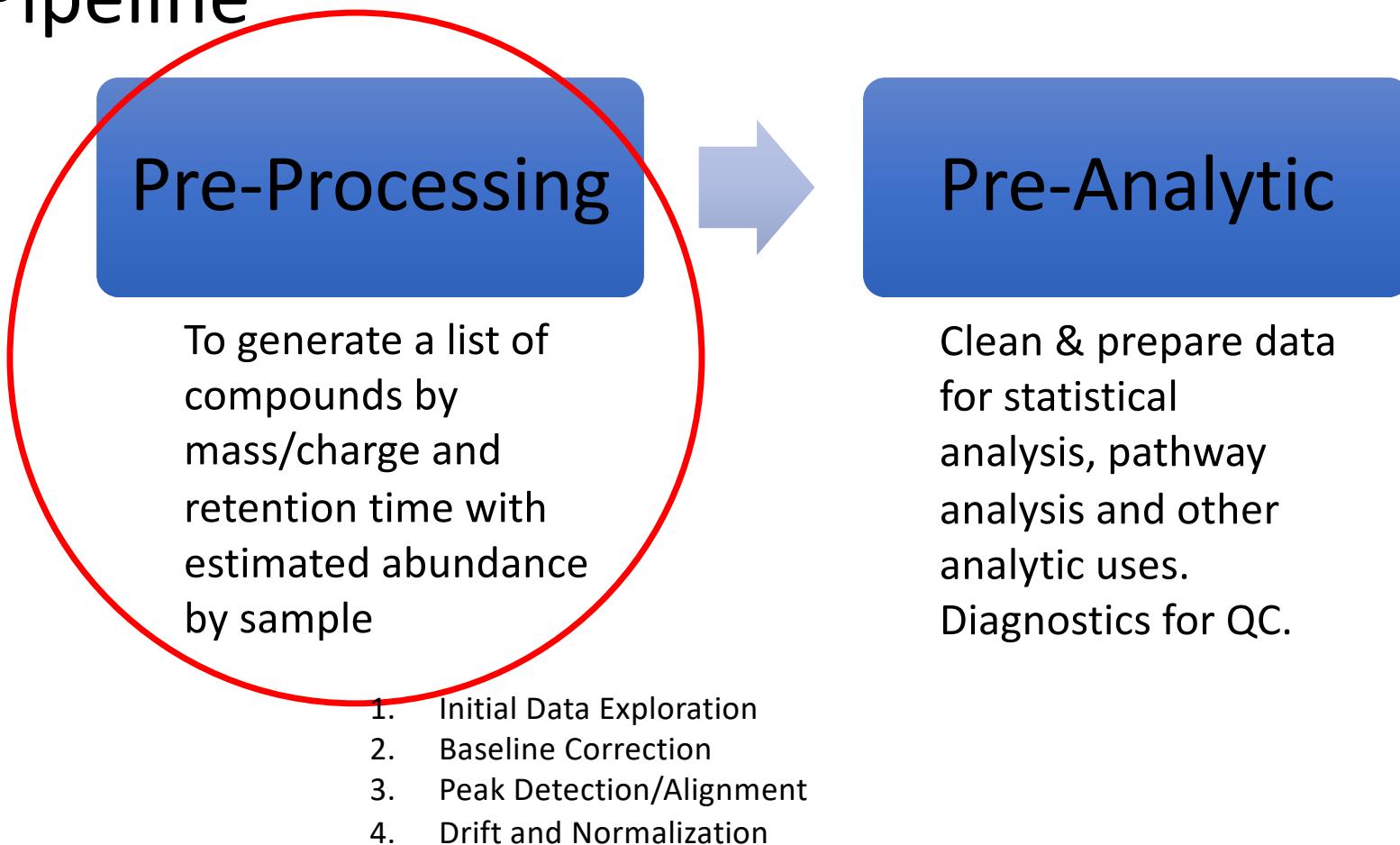


1. Pre-processing  
2. Pre-analytic

3. Identification

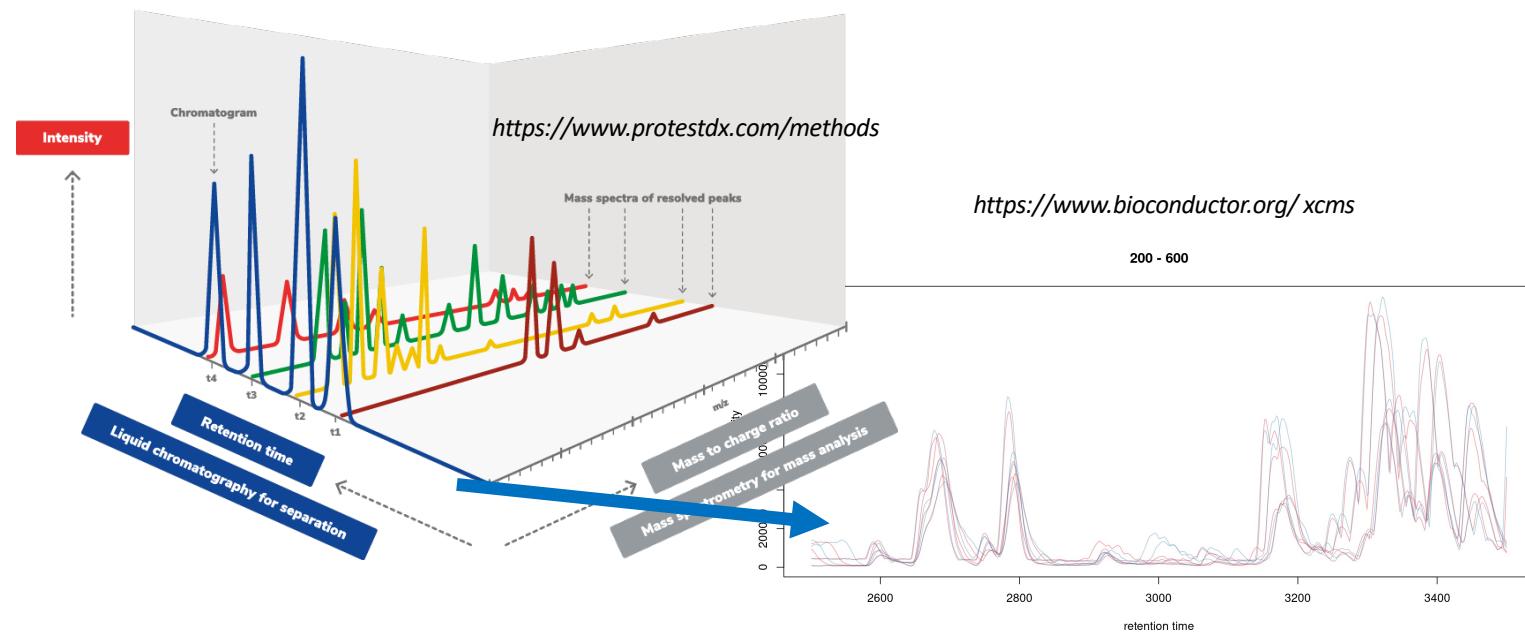
4. Downstream  
Analysis

# Pipeline



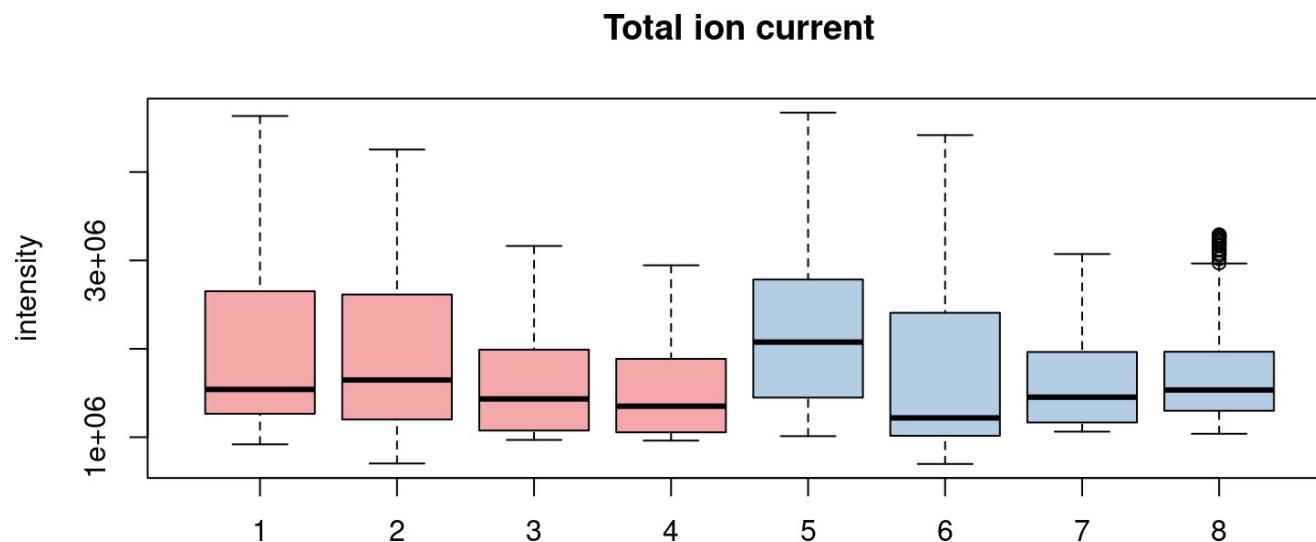
# 1. Initial Data Exploration

- To look at overall patterns & identify problematic samples
- Chromatogram building



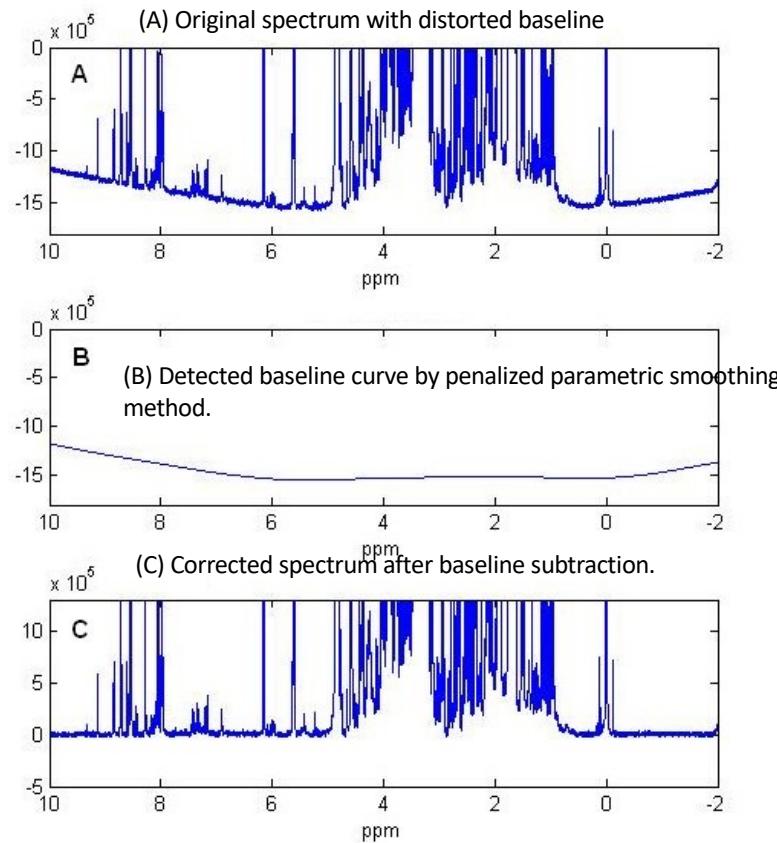
# 1. Initial Data Exploration

- Overall Intensity



<https://www.bioconductor.org/xcms>

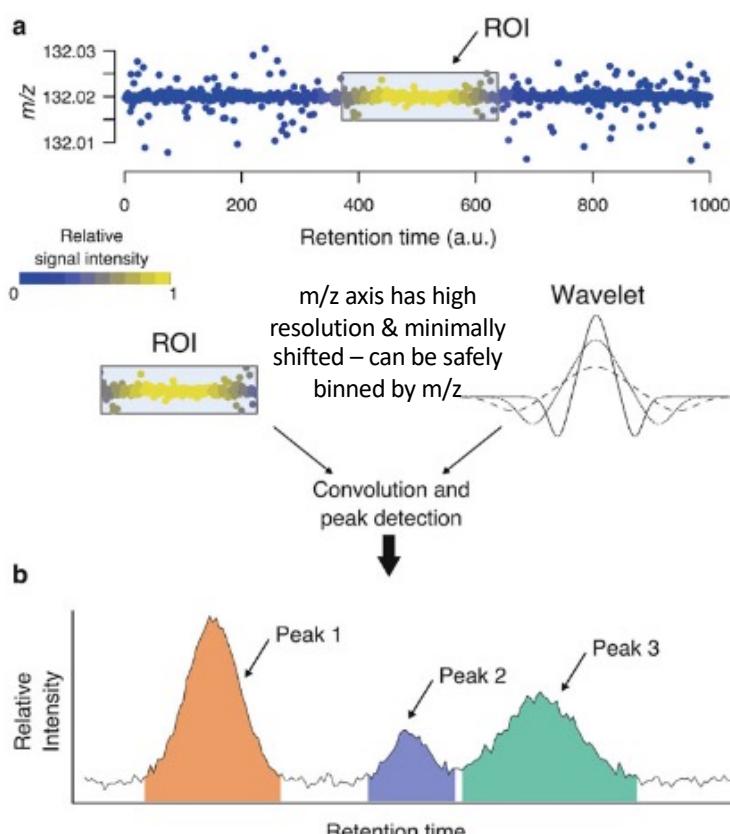
## 2. Baseline Correction



*Xi & Rocke (2008) BMC Bioinformatics*

# 3. Peak Detection

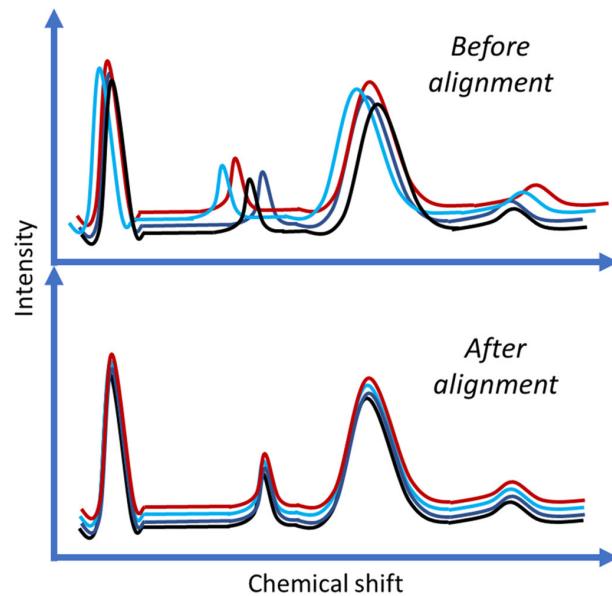
- By Sample
- Determine peaks above specified noise level that are true compounds
- Returns RT, m/z and integrated intensity for each peak



**Fig. 3** The centWave peak detection overview. First, regions of interest (ROI) are detected. For each detected ROI, a wavelet filter is applied and peaks within the ROI are detected

### 3. Peak Alignment

- Time of elution can vary across samples
- Alignment refers to find corresponding signals for the same compound across samples
- Often non-linear shifts

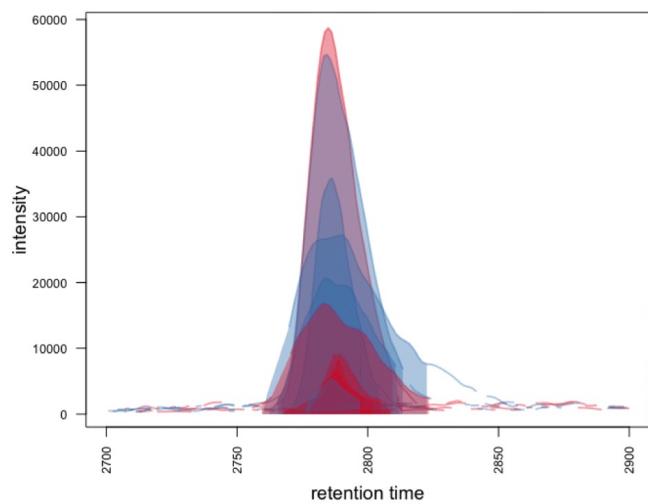


# 3. Peak Alignment

- Before peak detection – spectral alignment
  - Applying ‘warping’ methods (stretch/shrink), non-linear transformations to the RT to maximize correlation between spectra; Fast-Fourier transform methods
- After peak detection – peak alignment
  - Applied to peak coordinates, identify RT boundaries of peaks; use kernel density estimators

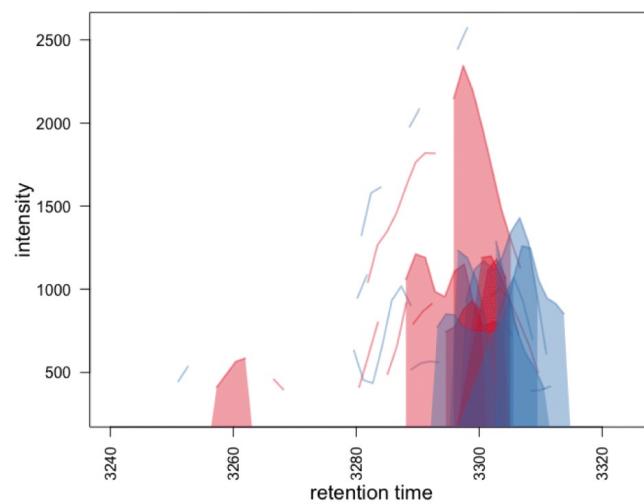
### 3. Peak Quality

Example of Good Peak Quality, m/z 335



(a) High Quality Feature

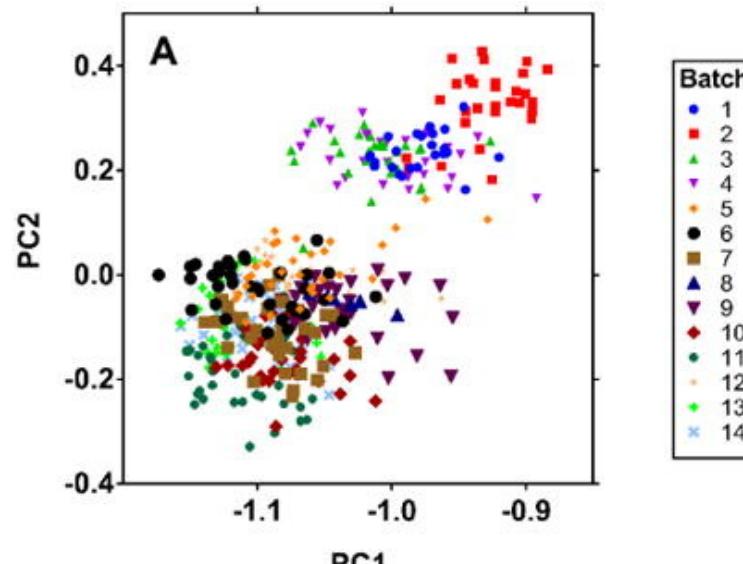
Example of Poor Peak Quality, m/z 246.6-246.8



(b) Low Quality Feature

## 4. Intensity drift

- Instrument response changes with time
- “Drift” can occur when samples are collected over multiple batches or long period of times



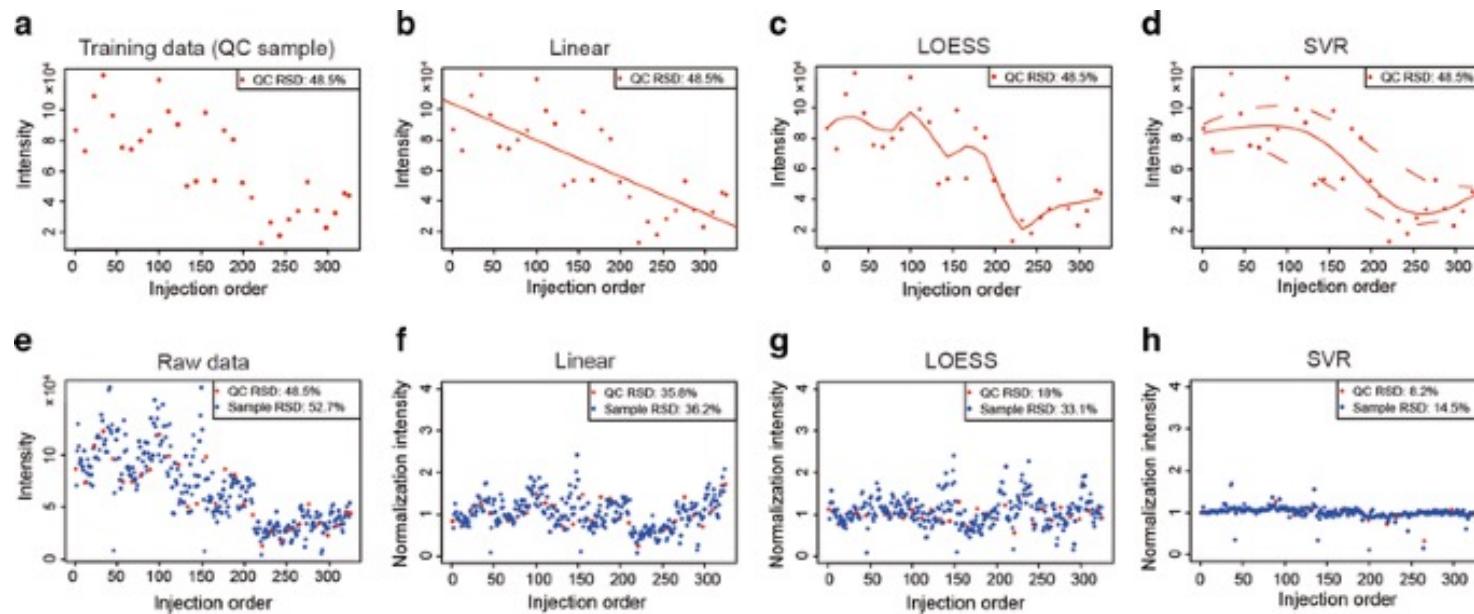
Thonusin et al., (2017) J. Chromatogr A

## 4. Strategies for Normalization

- Sum (or median) metabolite signal intensity normalization
  - assumes all metabolites experience same pattern of drift and that sum/median is approximately same in all samples
- Use Internal Standards (IS)
  - Not practical to include an isotope-labeled IS for all metabolites of interest.

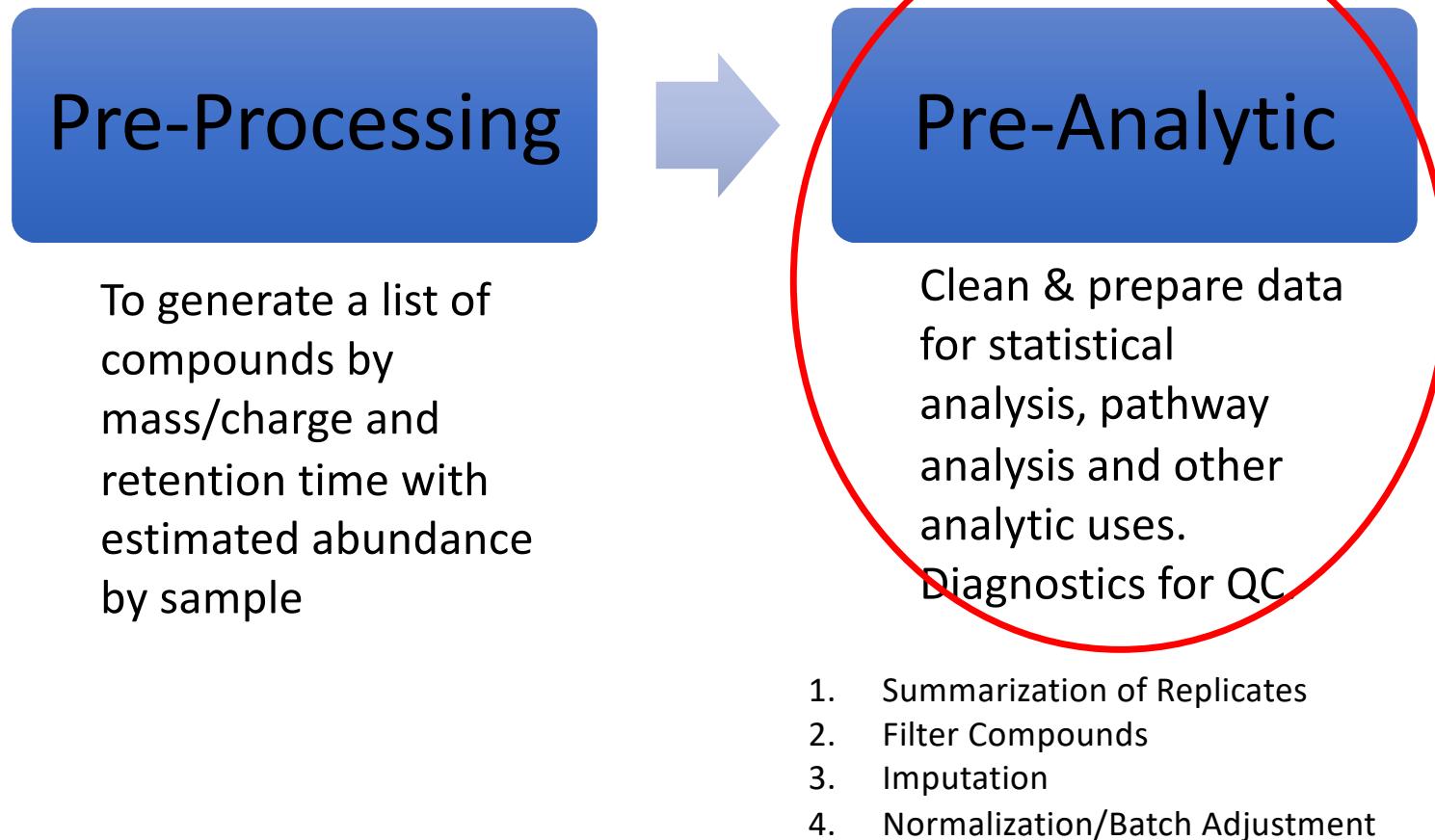
# 4. Strategies for Normalization

- Quality Control Samples
  - Included in each batch, can monitor and adjust for drift

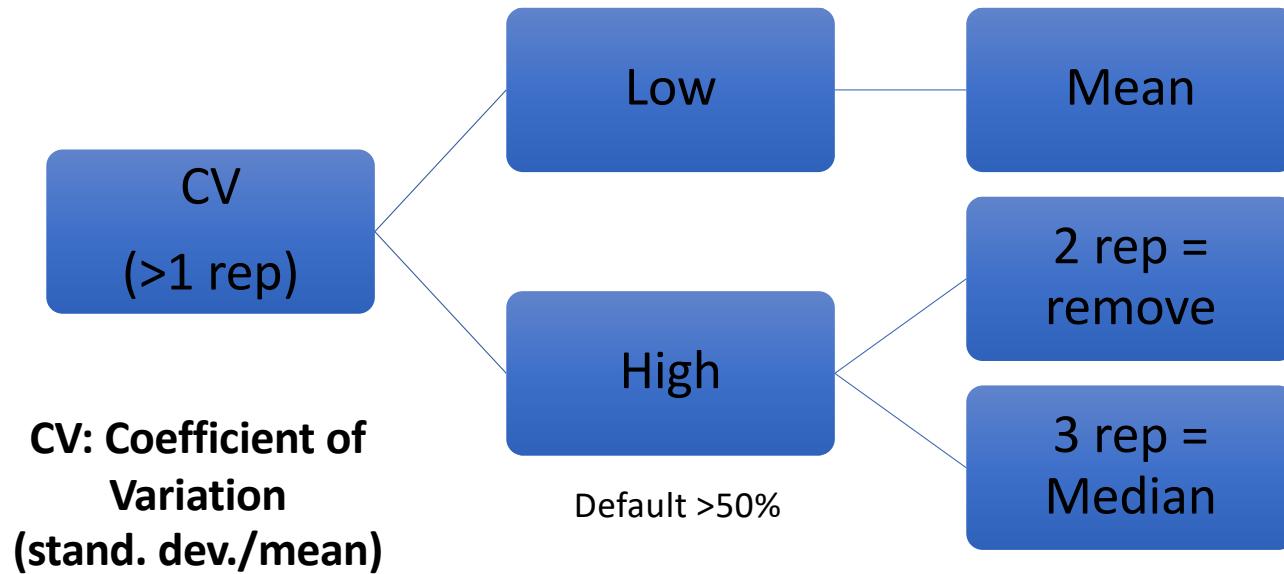


Shen et al. (2016) Metabolites

# Pipeline



# 1. Summarization of Replicates



Summarize technical replicates to result in one measure / subject / compound

*Hughes et al. (2014) Bioinformatics MSPrep*

# 1. Summarization of Replicates

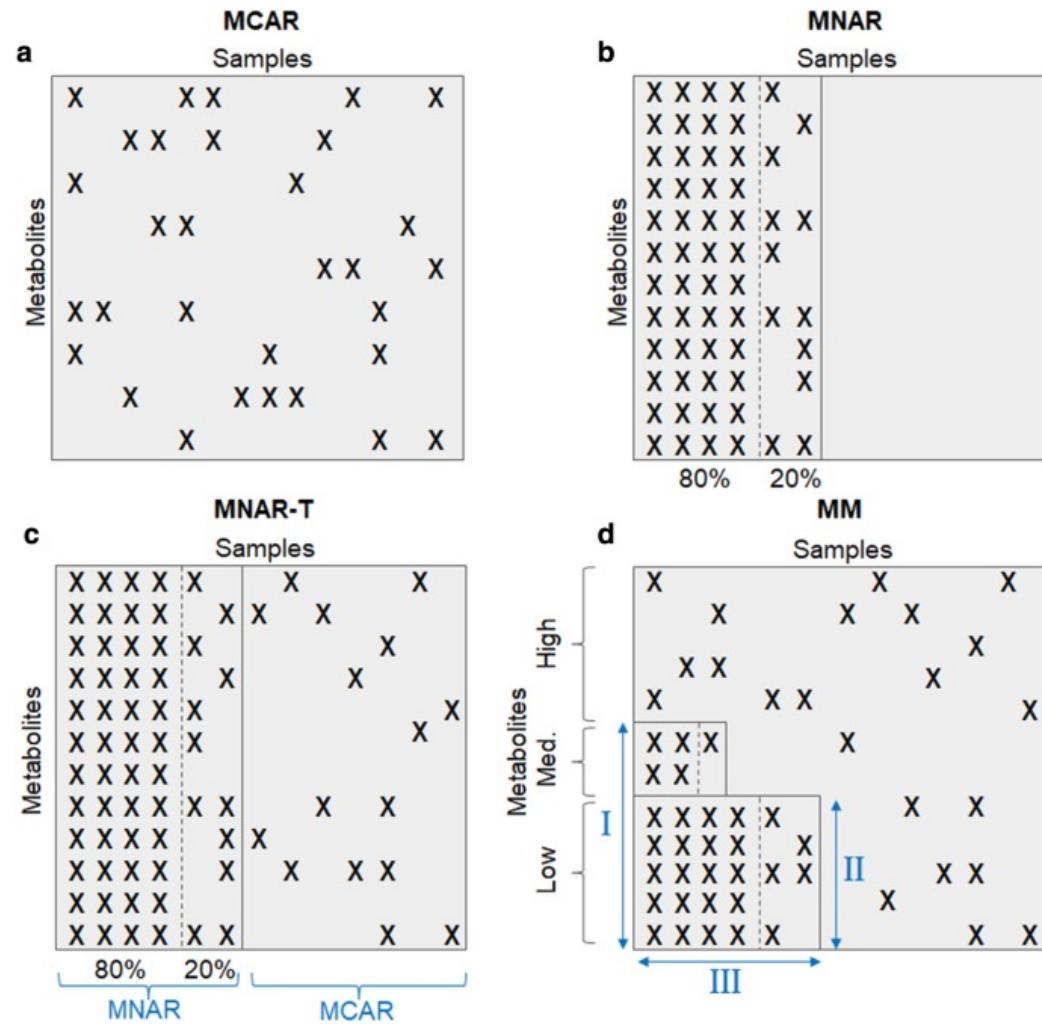
Most compounds with high CV have two consistent and one extreme observation.

Replicates

<i>Low</i>	<i>Med</i>	<i>High</i>	<i>Mean</i>	<i>SD</i>	<i>CV</i>	<i>Median</i>
<b>20839</b>	<b>30624</b>	251768	101077	130594	129.20%	30624
<b>101193</b>	<b>101761</b>	448747	217234	200497	92.30%	101761
9734	<b>263033</b>	<b>264957</b>	179241	146801	81.90%	263033
98592	<b>467745</b>	<b>468810</b>	345049	213439	61.86%	467745

## 2. Missing Data

- Three modes of missing data
  - Compound not present in sample
  - Compound present, but below detectable limit or specified noise level
  - Compound failed to be identified
    - Mechanical error (lab prep, LC column, coelution, etc.)
    - Errors in pre-processing (misalignment, retention time shift, mass accuracy, etc.)
- Missing Completely at Random (MCAR)
- Missing Not at Random (MNAR)



*Lee & Styczynski, (2018) Metabolomics*

## 2. Missing Data: Filtering Compounds

Goal: Remove compounds with too many missing values.

- Include compounds found in > X% of subjects (e.g., 80%)
  - Data driven approaches (Schiffman et al., (2019) BMC Bioinformatics)
- Removes noise detected as compounds, uncommon food or drug metabolites
- Simplifies downstream analysis
- Impact on downstream analysis
  - Too stringent reduces power and may miss important compounds
  - Too lenient increases false positives and drug/food metabolites, etc.

# 3. Imputation

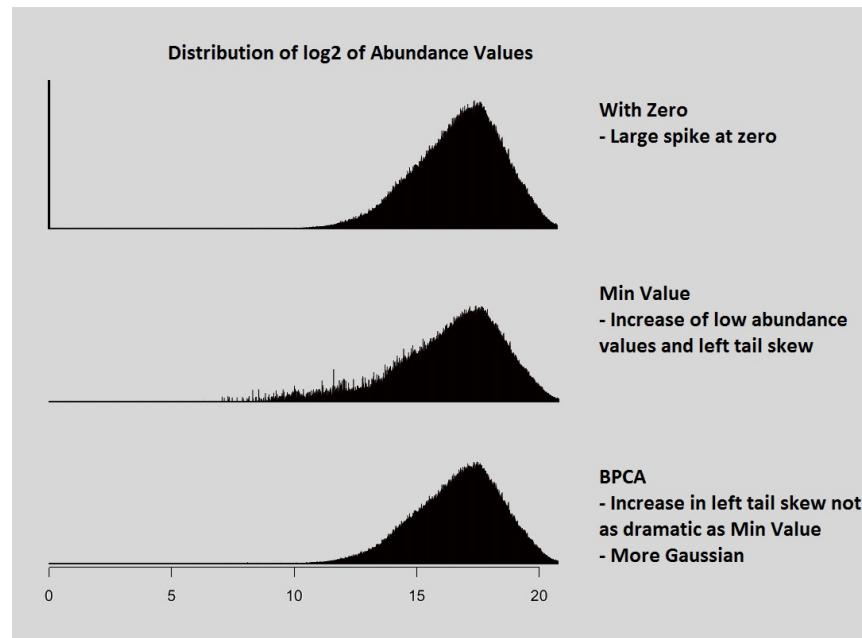
- Downstream analyses may need complete data (e.g., machine learning methods)
- Strategies for Imputation
  - I. With Zero – Zeros remain in dataset.
    - Normalization & some statistical analyses no longer valid.
    - Present/absent modeling.
  - II. Minimum Value – All missing values are replaced with  $\frac{1}{2}$  the minimum observed value for that compound.

# 3. Imputation

III. Random Forests

IV. K-Nearest  
Neighbors

V. Bayesian Principal  
Component  
Analysis (BPCA)



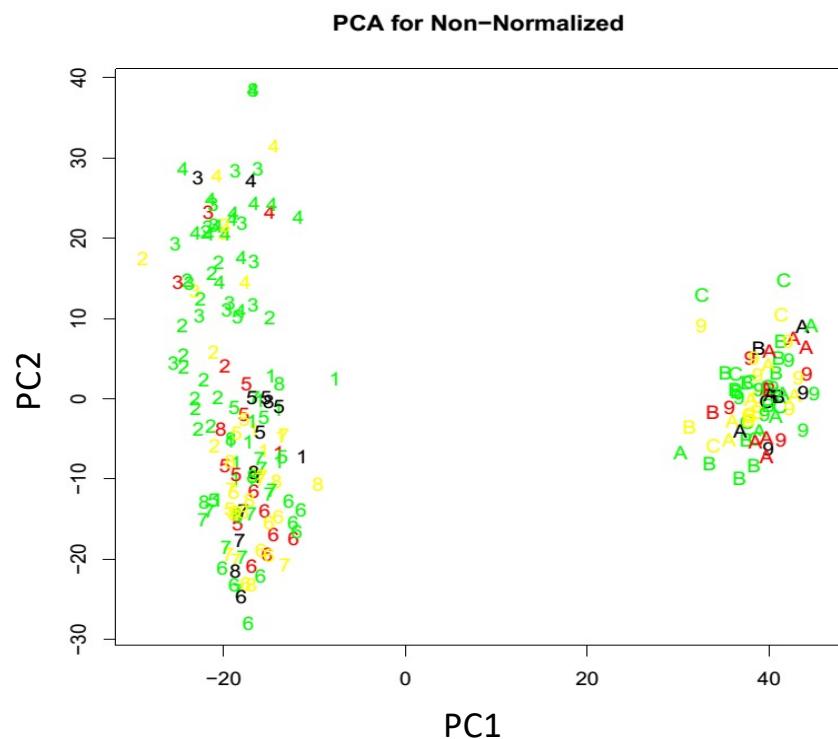
## 4. Normalization/Batch Effects

- Although QC samples can adjust for drift (see earlier slides), may have additional variation due to measured (e.g., batch effects) and unmeasured factors
- Utilize internal controls to measure variation.
  - Spiked controls or biological controls
- Data driven strategies
- Batch Effects
  - E.g., ComBat: Empirical Bayes Method to directly estimate and remove batch effect, can be applied post-hoc to any adjusted data

## 4. Normalization: Diagnostic Plots

- Principal Component Analysis (PCA)
- Boxplots of abundance by subject

Example: Color = Disease stage,  
Symbol = Batch



# 4. Normalization

## Strategies for Normalization

### A. No Normalization

- Data prepared by one technician, run in one batch, and/or specific compounds (small count) potentially may not need normalization

### B. Median Normalization

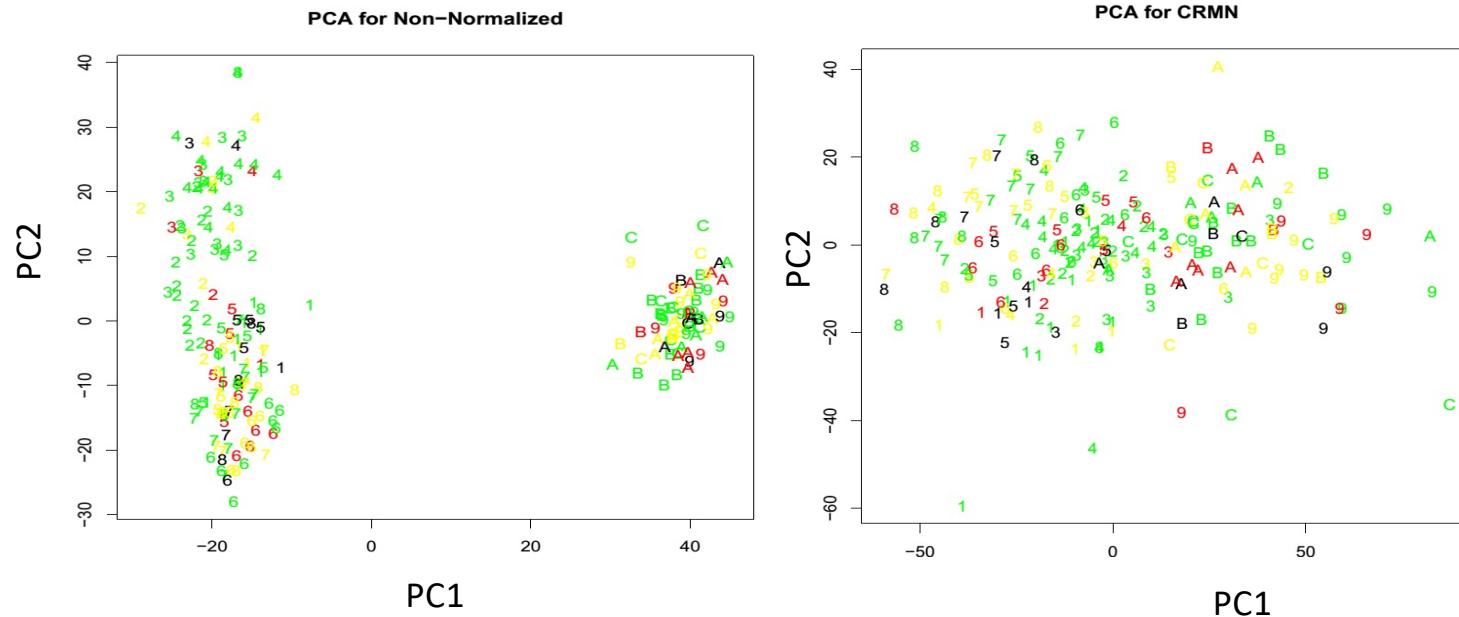
- Simplistic, but commonly used method
- Location shift accomplished by centering each subjects abundance on the median
- Q75 often selected instead of median if raw data is skewed or large quantity missing data

## 4. Normalization

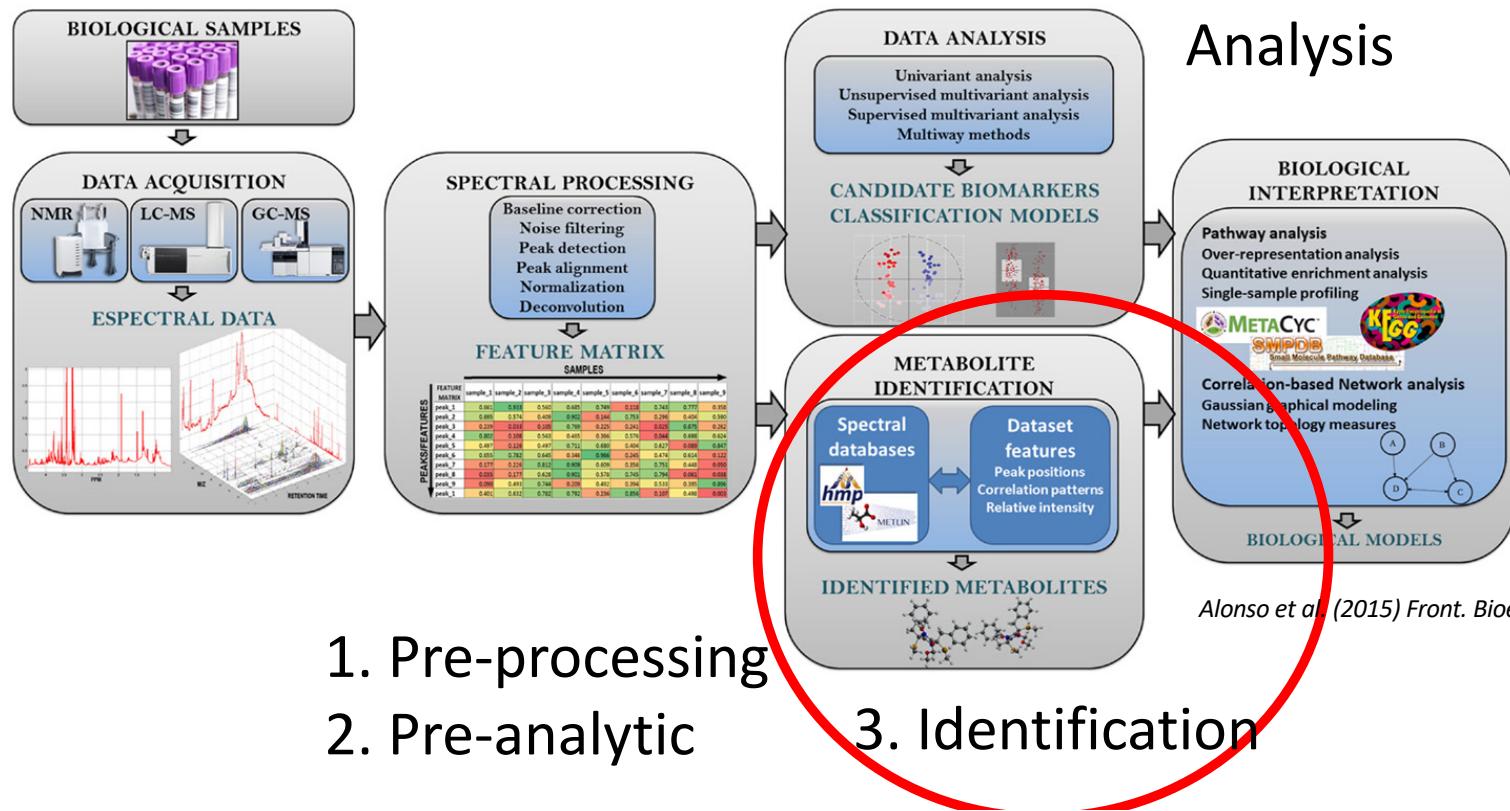
- C. Quantile Normalization
- D. Surrogate Variable Analysis (SVA)
- E. Remove Unwanted Variation – metabolomics (RUVm)
- G. Cross Contribution Normalization (CRMN)
  - Changing abundance of compounds is known to have potential to effect on other compounds
  - So some variation in internal standards could be due to variation in other compounds
  - Estimates the variance due to unmeasured sources similar to RUV, and estimates any significant cross contributing variance that comes from non-spiked compounds

# Example

Color = Disease Stage, Symbol = Batch



# Pipeline



# Identification

- **Goal:** identify most likely compound appearing at particular mass & retention time
- Anytime prior to pathway analysis
- Most based strictly off mass detected
- Often multiple potential matches and scoring algorithm can determine best match
- Major challenge in untargeted metabolomics

# Databases

- Facilitate metabolite identification
  - metabolite information - chemical formula, names and synonyms, physical and chemical properties
  - classification systems
  - structures and images
  - links to other databases (e.g., KEGG, PDB, ChEBI)
  - experimental data, mass spectrometry data

# Identification - Databases

- **HMDB** – Human Metabolome Database
  - <http://www.hmdb.ca/>
- **Metlin**
  - <https://metlin.scripps.edu/>
- **mzCloud**
  - <https://www.mzcloud.org>
- **MassBank/MoNA**
  - <http://www.massbank.jp>
- **NIST 17**
  - <https://chemdata.nist.gov>
- **Lipid MAPS** – Lipid Metabolites And Pathways Strategy
  - <http://www.lipidmaps.org/>

Database	Compounds with spectra	Number of Spectra
NIST 17	13,808	574,826
METLIN	>200,000	NA
MoNA	~75,000	261,917
mzCloud	> 8000	~2,000,000 <sup>a</sup>
HMDB	2265	22,247 <sup>b</sup>
LipidBlast	119,200	212,516

<sup>a</sup>The number of recalibrated spectra in mzCloud

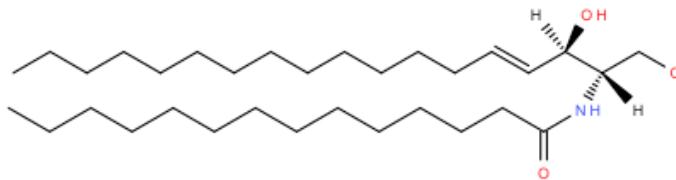
<sup>b</sup>Total experimental LC-MS/MS spectra in HMDB

# Example

Ceramide  
Cer(d18:1/14:0)  
from LIPID Maps

Structure database (LMSD)

Download file | MDLMOL MDLMOL files can be opened in various drawing programs. In ChemDr open as filetype "MDL Molfile".

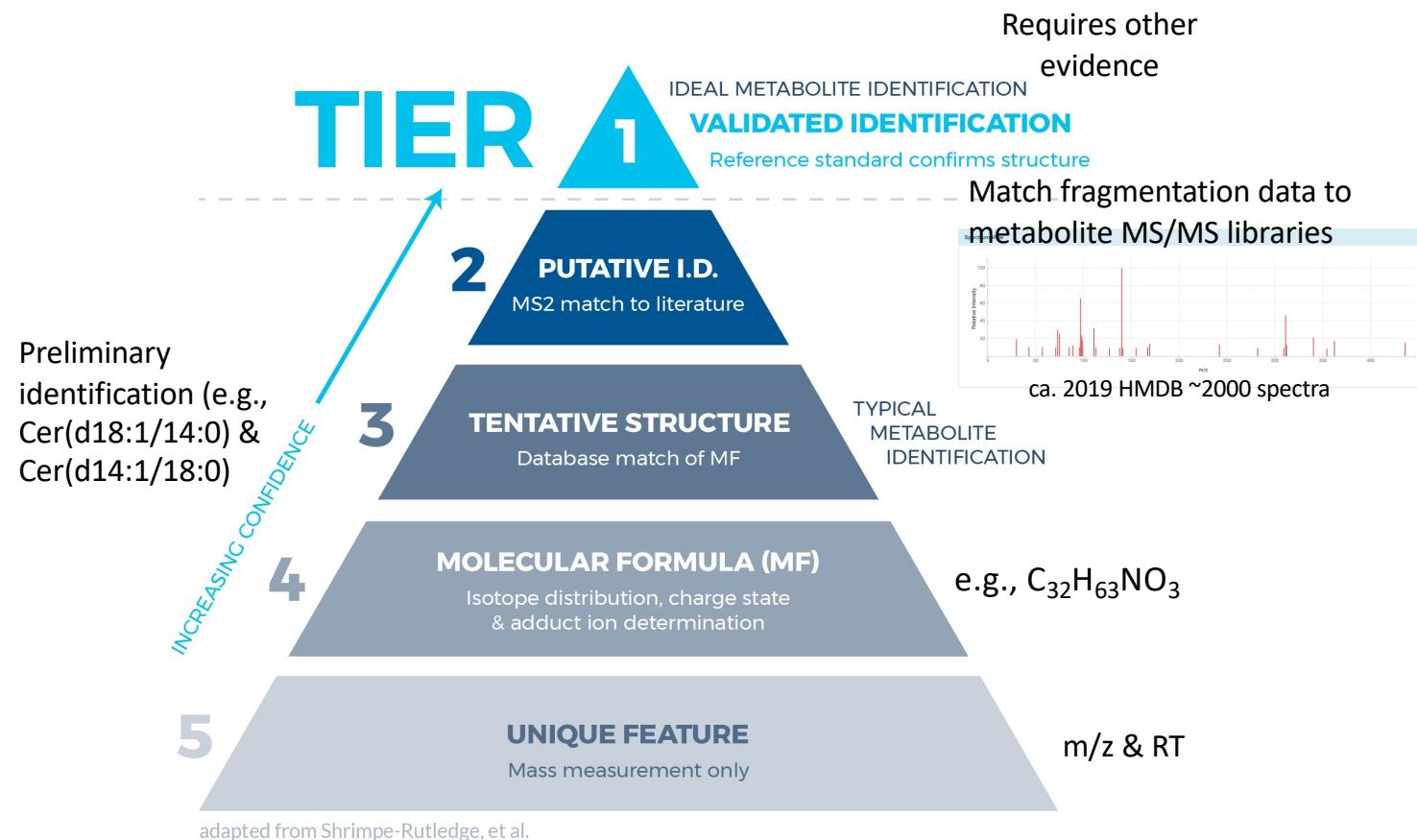


LM ID	LMSP02010001
Common Name	Cer(d18:1/14:0)
Systematic Name	N-(tetradecanoyl)-sphing-4-enine
Synonyms	C14 Cer; N-(tetradecanoyl)-ceramide; N-(myristoyl)-ceramide; Cer[NS]
Exact Mass	509.4808 (neutral) Calculate m/z: (Select m/z)
Formula	C <sub>32</sub> H <sub>63</sub> NO <sub>3</sub>
Category	Sphingolipids [SP]
Main Class	<a href="#">Ceramides [SP02]</a>
Sub Class	<a href="#">N-acylsphingosines (ceramides) [SP0201]</a>
PubChem Compound ID (CID)	<a href="#">5282310</a>
METABOLOMICS ID	-
KEGG ID	-
HMDB ID	<a href="#">HMDB11773</a>

<http://www.lipidmaps.org/data/LMSDRecord.php?LMID=LMSP02010001>



# Levels of Identification



lipidmaps.org/data/structure/LMSDSearch.php

Apps Google CU Denver Access Home Page GitHub Canvas PubMed Yahoo Mail Weather mil

LIPID MAPS® Lipidomics Gateway Home Updates ▾ Reso

# Structure Database (LMSD)

## Text/Ontology-based search

LM ID

Name (Common, Systematic, or Synonym)

Mass  +/- 0.5

Formula

Category

Main class

Sub class

Ontology search parameters

Group:

Group:

Group:

Group:

Include  All records  Curated records only  Computationally generated records only

Records per page

Sort by

LMSD: Text/Ontology-based search results

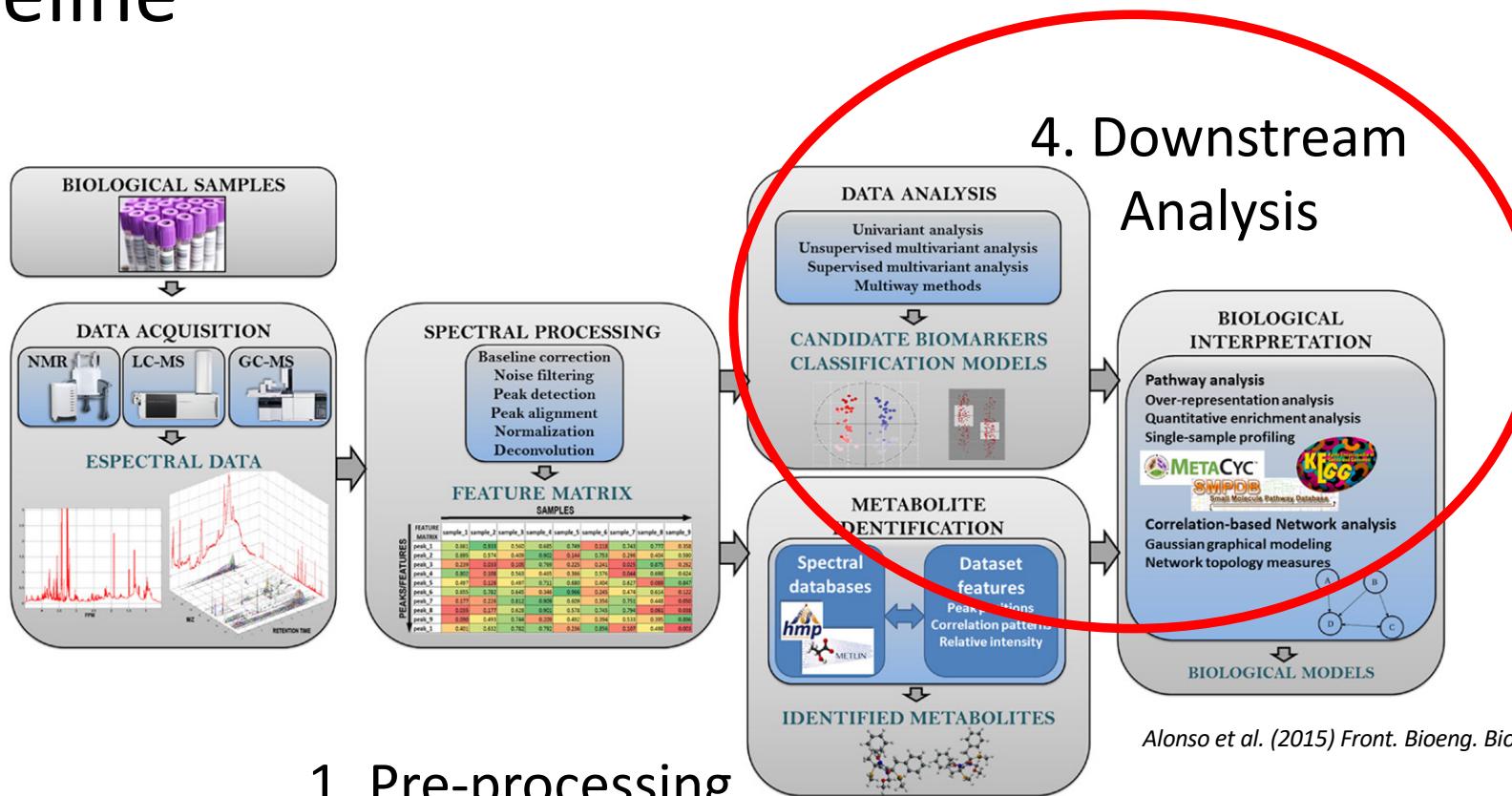
[Modify Search](#)

LM_ID	Common Name	Systematic Name	Species Shorthand	Formula	Mass	Main Class	Sub Class
LMFA01030488	3-heptynoic acid	3-heptynoic acid	-	C <sub>7</sub> H <sub>10</sub> O <sub>2</sub>	126.0681	Fatty Acids and Conjugates [FA01]	Unsaturated fatty acids [FA0103]
LMFA01030489	4-heptynoic acid	4-heptynoic acid	-	C <sub>7</sub> H <sub>10</sub> O <sub>2</sub>	126.0681	Fatty Acids and Conjugates [FA01]	Unsaturated fatty acids [FA0103]
LMFA01030490	6-heptynoic acid	6-heptynoic acid	-	C <sub>7</sub> H <sub>10</sub> O <sub>2</sub>	126.0681	Fatty Acids and Conjugates [FA01]	Unsaturated fatty acids [FA0103]
LMFA01030792	4-Methyl-3Z,5-hexadienoic acid	4-Methyl-3Z,5-hexadienoic acid	-	C <sub>7</sub> H <sub>10</sub> O <sub>2</sub>	126.0681	Fatty Acids and Conjugates [FA01]	Unsaturated fatty acids [FA0103]
LMFA05000114	2,4-Dimethyl-2E,4E-hexadien-1-ol	2,4-Dimethyl-2E,4E-hexadien-1-ol	-	C <sub>8</sub> H <sub>14</sub> O	126.1045	Fatty alcohols [FA05]	-
LMFA05000488	1,5E-Octadien-3-ol	1,5E-Octadien-3-ol	-	C <sub>8</sub> H <sub>14</sub> O	126.1045	Fatty alcohols [FA05]	-
LMFA05000493	1,5Z-Octadien-3-ol	1,5Z-Octadien-3-ol	-	C <sub>8</sub> H <sub>14</sub> O	126.1045	Fatty alcohols [FA05]	-
LMFA06000027	2-heptenedial	2-heptenedial	-	C <sub>7</sub> H <sub>10</sub> O <sub>2</sub>	126.0681	Fatty aldehydes [FA06]	-
LMFA06000029	2-octenal	2-octenal	-	C <sub>8</sub> H <sub>14</sub> O	126.1045	Fatty aldehydes [FA06]	-
LMFA06000030	3-octenal	3-octenal	-	C <sub>8</sub> H <sub>14</sub> O	126.1045	Fatty aldehydes [FA06]	-
LMFA06000031	4-octenal	4-octenal	-	C <sub>8</sub> H <sub>14</sub> O	126.1045	Fatty aldehydes [FA06]	-
LMFA06000032	5-octenal	5-octenal	-	C <sub>8</sub> H <sub>14</sub> O	126.1045	Fatty aldehydes [FA06]	-
LMFA06000033	6-octenal	6-octenal	-	C <sub>8</sub> H <sub>14</sub> O	126.1045	Fatty aldehydes [FA06]	-
LMFA07010967	Methyl sorbate	methyl (2E,4E)-hexa-2,4-dienoate	-	C <sub>7</sub> H <sub>10</sub> O <sub>2</sub>	126.0681	Fatty esters [FA07]	Wax monoesters [FA0701]
LMFA11000323	1-Nonene	1-Nonene	-	C <sub>9</sub> H <sub>18</sub>	126.1409	Hydrocarbons [FA11]	-
LMFA11000594	3-ethyl, 4,4-dimethyl-2-Pentene	3-ethyl, 4,4-dimethyl-2-Pentene	-	C <sub>9</sub> H <sub>18</sub>	126.1409	Hydrocarbons [FA11]	-
LMFA11000629	Isopropyl-cyclohexane	1-methylethyl-cyclohexane	-	C <sub>9</sub> H <sub>18</sub>	126.1409	Hydrocarbons [FA11]	-
LMFA11000632	1,1,2-trimethylcyclohexane	1,1,2-trimethylcyclohexane	-	C <sub>10</sub> H <sub>20</sub>	126.1409	Hydrocarbons [FA11]	-

Search using  
 $C_{16}H_{26}O_3$

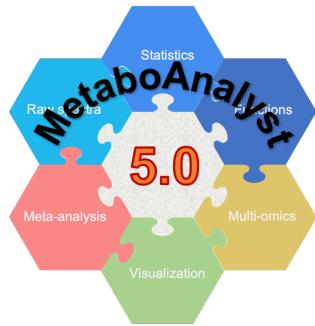
LM_ID	Common Name	Systematic Name	Species Shorthand	Formula
LMFA01020279	10,11-epoxy-3,7,11-trimethyl-2E,6E-tridecadienoic acid	10,11-epoxy-3,7,11-trimethyl-2E,6E-tridecadienoic acid	-	$C_{16}H_{26}O_3$
LMFA01031023	7-Hydroxy-10E-Hexadecen-8-ynoic acid	7-Hydroxy-10E-Hexadecen-8-ynoic acid	-	$C_{16}H_{26}O_3$
LMFA01050143	Tetranor-12R-HETE	8R-hydroxy-4Z,6E,10Z-hexadecatrienoic acid	-	$C_{16}H_{26}O_3$
LMFA01050562	9S-HHTrE	9S-hydroxy-6Z,10E,12Z-hexadecatrienoic acid	-	$C_{16}H_{26}O_3$
LMFA01060229	-	4-oxo-2E,15-hexadecenoic acid	-	$C_{16}H_{26}O_3$
LMFA01070019	(2E,6E,10R,11S)-10,11-epoxy-3,7,11-trimethyltrideca-2,6-dienoic acid	10R,11S-epoxy-3,7,11-trimethyl-2E,6E-tridecadienoic acid	-	$C_{16}H_{26}O_3$
LMFA01070023	13,14-epoxy-7Z,10Z-Hexadecenoic acid	13,14-epoxy-7Z,10Z-Hexadecenoic acid	-	$C_{16}H_{26}O_3$
LMFA01150024	5-(3,4-dimethyl-5-pentylfuran-2-yl)-pentanoic acid	5-(5-pentyl-3,4-dimethylfuran-2-yl)-pentanoic acid	-	$C_{16}H_{26}O_3$
LMFA02010008	(9R,13R)-1a,1b-dinor-10,11-dihydro-12-oxo-15-phytoenoic acid	(1R,2R)-3-oxo-2-(2'Z-pentenyl)cyclopentanehexanoic acid	-	$C_{16}H_{26}O_3$
LMFA02010009	(9S,13S)-1a,1b-dinor-10,11-dihydro-12-oxo-15-phytoenoic acid	(1S,2S)-3-oxo-2-(2'Z-pentenyl)cyclopentanehexanoic acid	-	$C_{16}H_{26}O_3$
LMPR0103010004	Juvenile hormone III (W)	-	-	$C_{16}H_{26}O_3$

# Pipeline



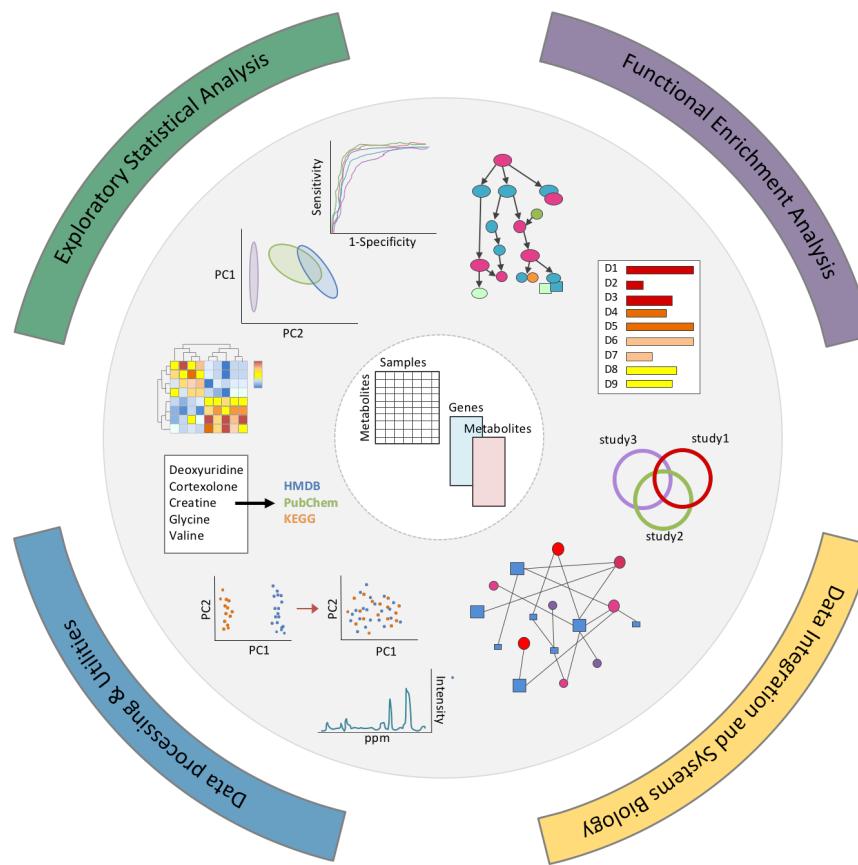
# Downstream Analysis

- Statistical analysis
  - t-test, regression, ANOVA
  - multiple testing correction (e.g., Bonferroni, false discovery rate)
- Machine learning
  - dimension reduction (principal component analysis), hierarchical clustering
  - logistic regression, random forests, support vector machines, neural networks, deep learning
  - cross-validation



# MetaboAnalyst 5.0

user-friendly, streamlined metabolomics data analysis



# Pathway Databases

- **KEGG** - Kyoto Encyclopedia of Genes and Genomes
  - <http://www.genome.jp/kegg/>
  - ~530 pathways
- **MetaCyc**
  - <https://MetaCyc.org>
  - ~2600 pathways
- **Reactome**
  - <https://reactome.org/>
  - ~20,000 pathways

# Enrichment Analysis

## ChemRICH

- Chemical Similarity Enrichment Analysis for Metabolomics

### Why use ChemRICH ?

Metabolomics aims to answer a fundamental question in biology: How does metabolism change under genetic, environmental or phenotypic perturbations?

Combining several metabolomics assays can yield datasets for more than 1,000 structurally identified metabolites per study. However, biological interpretations of metabolic regulation in these datasets is hindered by the limitations of current pathway definitions as well as inherent limits of pathway enrichment statistics.

ChemRICH, a statistical enrichment approach that is based on chemical similarity rather than sparse biochemical knowledge annotations. ChemRICH utilizes chemical ontologies and structural similarity to group metabolites. Unlike pathway mapping, this strategy yields study-specific, non-overlapping sets of all identified metabolites. Subsequent enrichment statistics is superior to pathway enrichments because ChemRICH sets have a self-contained size where p-values do not rely on the size of a background database. For more details - see [ChemRICH article](#)

### Input file structure

	A	B	C	D	E	F
	Compound Name	InChiKeys	Pubchem ID	SMILES	pvalue	foldchange
1	SM (442,2) B	DAOCOGIMBILZYOH-GXJPFUDISA-N	52931217	CCCCCCCCCCCCCCCCCCCCCCC=C(O)N[C]13.32E-13	0.30587795	
2	5M (442,2)	DACOGIMBILZYOH-GXJPFUDISA-N	44280128	CCCCCCCCCC=O</[C@H]([C@@H]1C3.71E-13	0.279551732	
3	N-Tetradecenoyl-4-spiro-2H-heptal-3-one-O-	QYKFWS55SA-N	86354	CCCCCCCCCCCCC=C(O)OC(=O)[C@H]1C3.50E-12	0.36195589	
4	LPC (16:0)	ASWBNKHC2GQVIV-UHF7FAOYA-N	5497109	CCCCCCCCCCCCC=C(O)OC(C@H)COP(=O)([C@H]1C5.38E-12	0.369242722	
5	PC (18:1)/16:0]	WTJKGGKOPKQKULL-VYOB0KEKSAA-N				

Structure of the input file. See here an [example data file](#) you can use as template.

### The input file must have 6 columns, in this order:

- Column 1 = Compound Name
- Column 2 = InChiKeys
- Column 3= Pubchem ID
- Column 4 = SMILES
- Column 5 = pvalue
- Column 6= foldchange

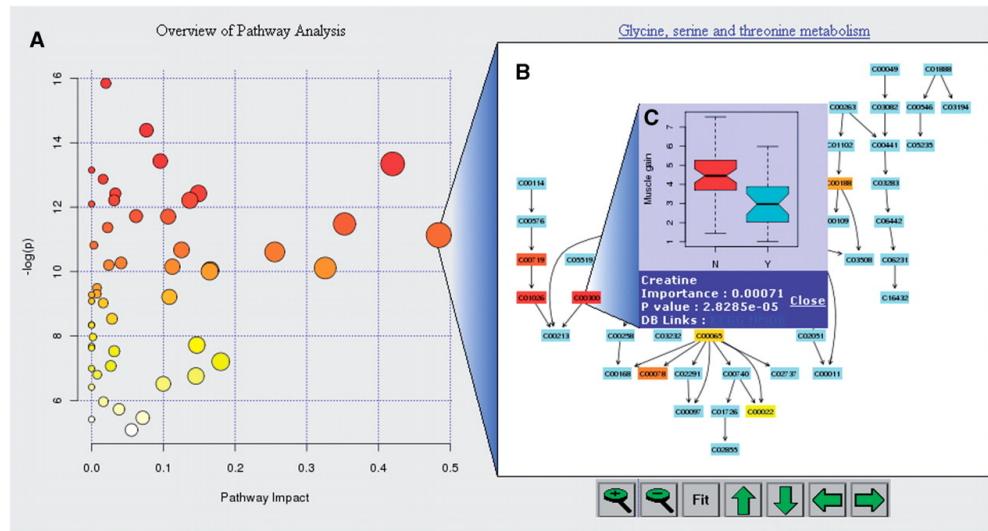
You can obtain PubChem CIDs and InChi keys from [the Chemical Translation Service](#). You can obtain SMILES from [the PubChem Identifier Exchanger tool](#).

**You can also use this excel file for getting identifiers and SMILES codes for your compounds. [ChemRICH Metabolite Identifier File](#) .**

**Input file must satisfy these conditions**

# Pathway Analysis

metPA in MetaboAnalyst



**Pathway Impact:**  
Examines compounds' position in pathway, centrality/hub



The screenshot shows the homepage of the Metabolomics Workbench. At the top left is the logo "Metabolomics Workbench". The main title "METABOLOMICS WORKBENCH" is prominently displayed in the center. On the right side, there are links for "Log in / Register" and a search bar with the placeholder "Search the Metabolomics Workbench". Below the header, a navigation bar includes links for Home, Data Repository, Databases, Protocols, Tools, Training / Events, About, and Search. A welcome message at the bottom states: "Welcome to the UCSD Metabolomics Workbench, a resource sponsored by the Common Fund of the National Institutes of Health."

### National Metabolomics Data Repository

[Upload and Manage Studies](#) | [Browse and Search Studies](#) | [Analyze Studies](#)

As of 04/08/21 a total of 1657 studies have been processed by the National Metabolomics Data Repository (NMDR). There are 1396 publicly available studies and the remainder (261) will be made available subject to their embargo dates.

#### Recently released studies on NMDR

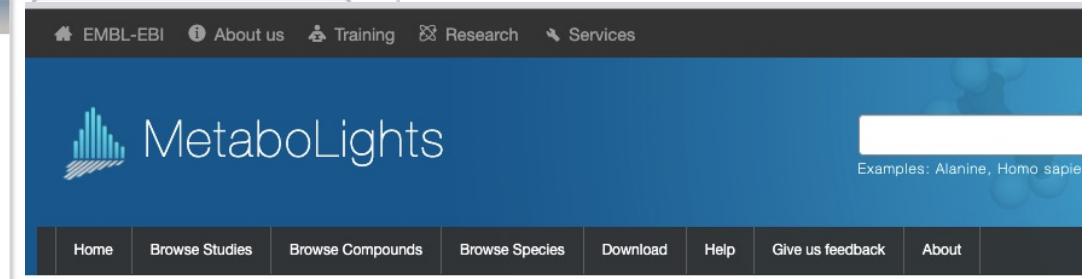
**ST001260** - Metabolic changes of *Fusobacterium nucleatum* when co-cultured with other oral microbes (part-I); *Fusobacterium nucleatum*; [Osaka University](#)

**ST001519** - Stool metabolites of known identity profiled using hybrid nontargeted methods (part-II); *Homo sapiens*; [Broad Institute of MIT and Harvard](#)

**ST001520** - Stool unknowns profiled using hybrid nontargeted methods (part-II); *Homo sapiens*; [Broad Institute of MIT and Harvard](#)

**ST001521** - Plasma metabolites of known identity profiled using hybrid nontargeted methods (part-III); *Homo sapiens*; [Broad Institute of MIT and Harvard](#)

<https://www.metabolomicsworkbench.org/>



The screenshot shows the homepage of MetaboLights. At the top, there is a navigation bar with links for EMBL-EBI, About us, Training, Research, and Services. The main title "MetaboLights" is centered on a blue background. Below the title is a search bar with the placeholder "Examples: Alanine, Homo sapiens". A secondary navigation bar at the bottom includes links for Home, Browse Studies, Browse Compounds, Browse Species, Download, Help, Give us feedback, and About.

## MetaboLights

MetaboLights is a database for Metabolomics experiments and derived information. The database is cross-species, cross-technique and covers metabolite structures and their reference spectra as well as their biological roles, locations and concentrations, and experimental data from metabolic experiments. MetaboLights is the recommended Metabolomics repository for a number of leading journals.

[More about us](#)

[Quick tour](#)

  
**Study**  
[BROWSE](#)

  
**Compound Library**  
[COMPOUNDS](#)

  
**Training**  
[METABOLOMICS TRAIN ONLINE](#)

<https://mscat.metabolomicsworkbench.org/>

Tool Name	Website	Publication	Last Updated	Operating System	User Interface	Molecule Type
filter data...						
CluMSID	<a href="#">Link</a>	<a href="#">Link</a>	2020-10-16	Windows, Mac, Linux	CLI	metabolomics
GNPS-MassIVE	<a href="#">Link</a>	<a href="#">Link</a>	2020-10-12	Windows, Mac, Linux	web	metabolomics
mineXpert	<a href="#">Link</a>	<a href="#">Link</a>	2020-10-12	Windows, Mac, Linux	GUI	metabolomics, protec
nPYc-Toolbox	<a href="#">Link</a>	<a href="#">Link</a>	2020-10-07	Mac, Linux, Windows	CLI	metabolomics
EI-MAVEN	<a href="#">Link</a>	<a href="#">Link</a>	2020-09-28	Windows, Mac	GUI	metabolomics
VSClust	<a href="#">Link</a>	<a href="#">Link</a>	2020-08-12	Windows, Mac, Linux	web	metabolomics, genon
mixOmics	<a href="#">Link</a>	<a href="#">Link</a>	2020-08-09			genomics, transcriptc

# “Metabolomics: the Superglue of Omics”

Giera, Spilker, Siuzdak (2018)

