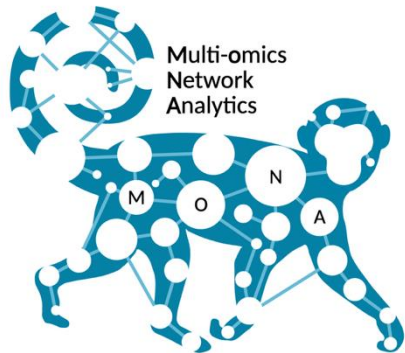


# Informatics Platform and MoNA

An introduction to Metabolomics

Date: 03 June 2026



Henry Webel, Senior Data Scientist

Felicia Cara Schulz, PhD student

Maria Barranco Altirriba, PostDoc

Alberto Pallejà, Team Lead Data Science Platform

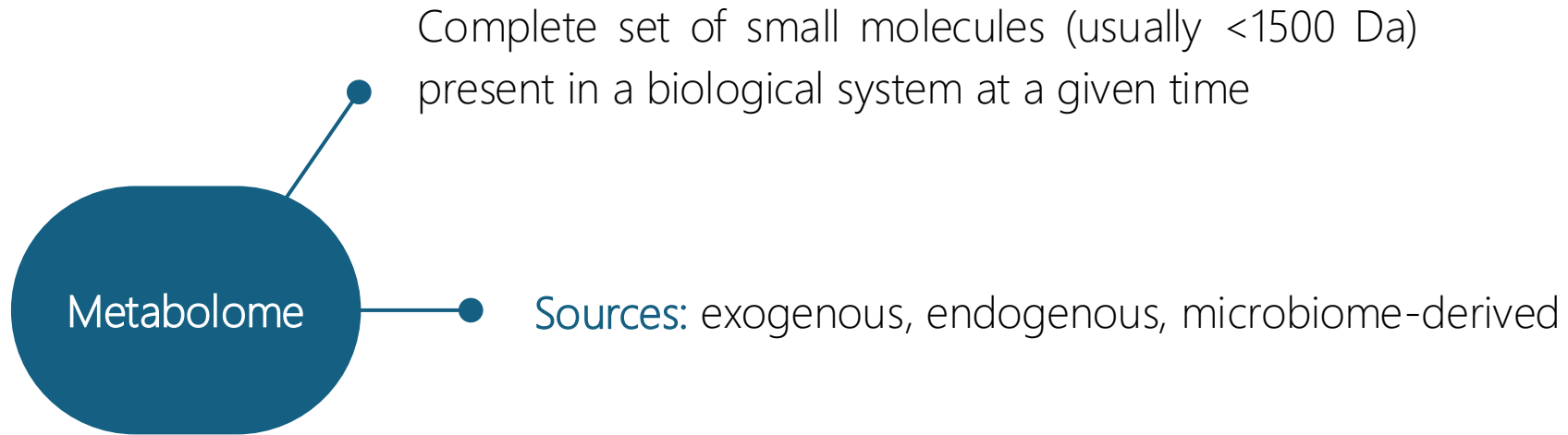
# Introduction to metabolomics

Complete set of small molecules (usually <1500 Da)  
present in a biological system at a given time

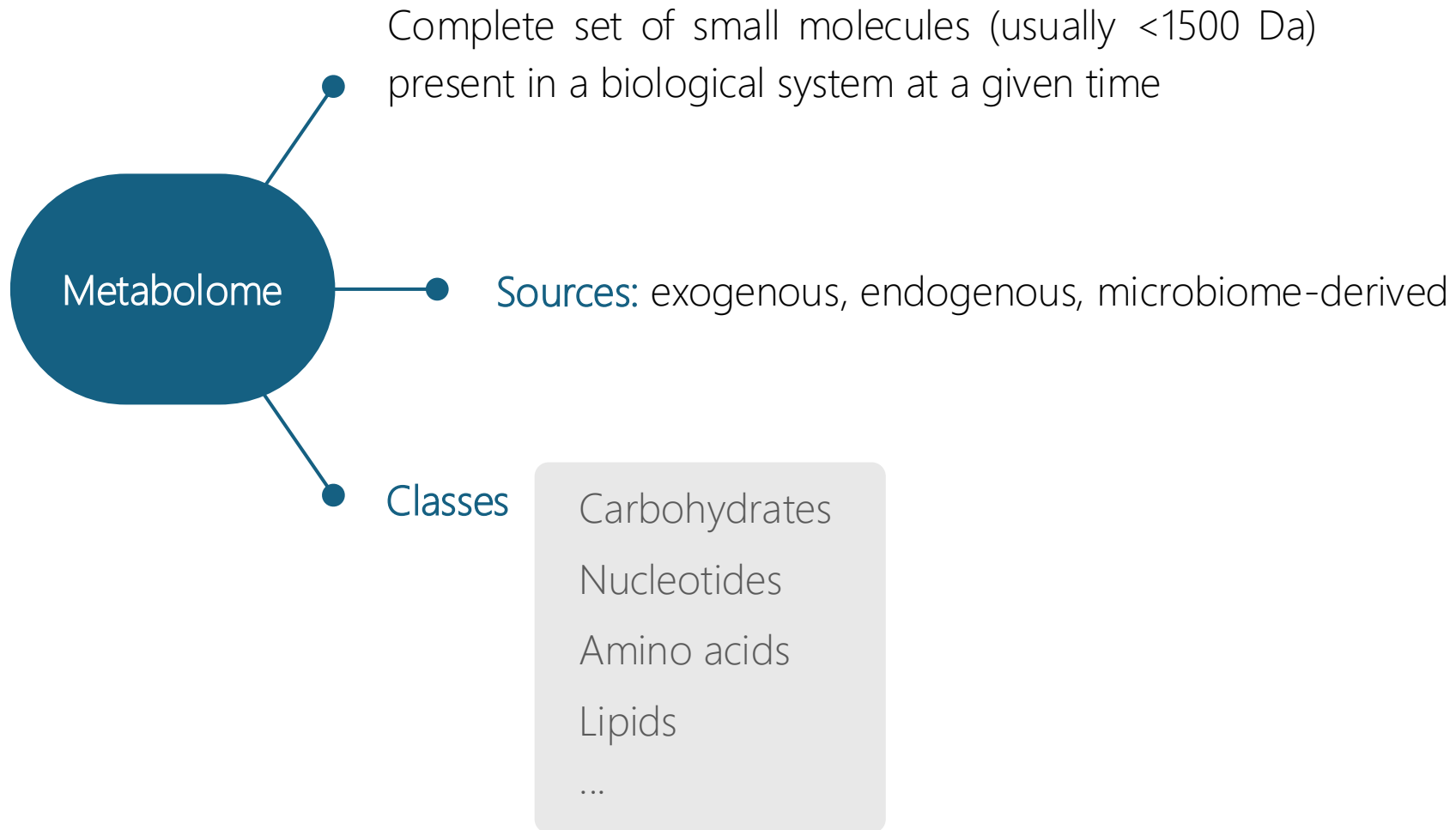


Metabolome

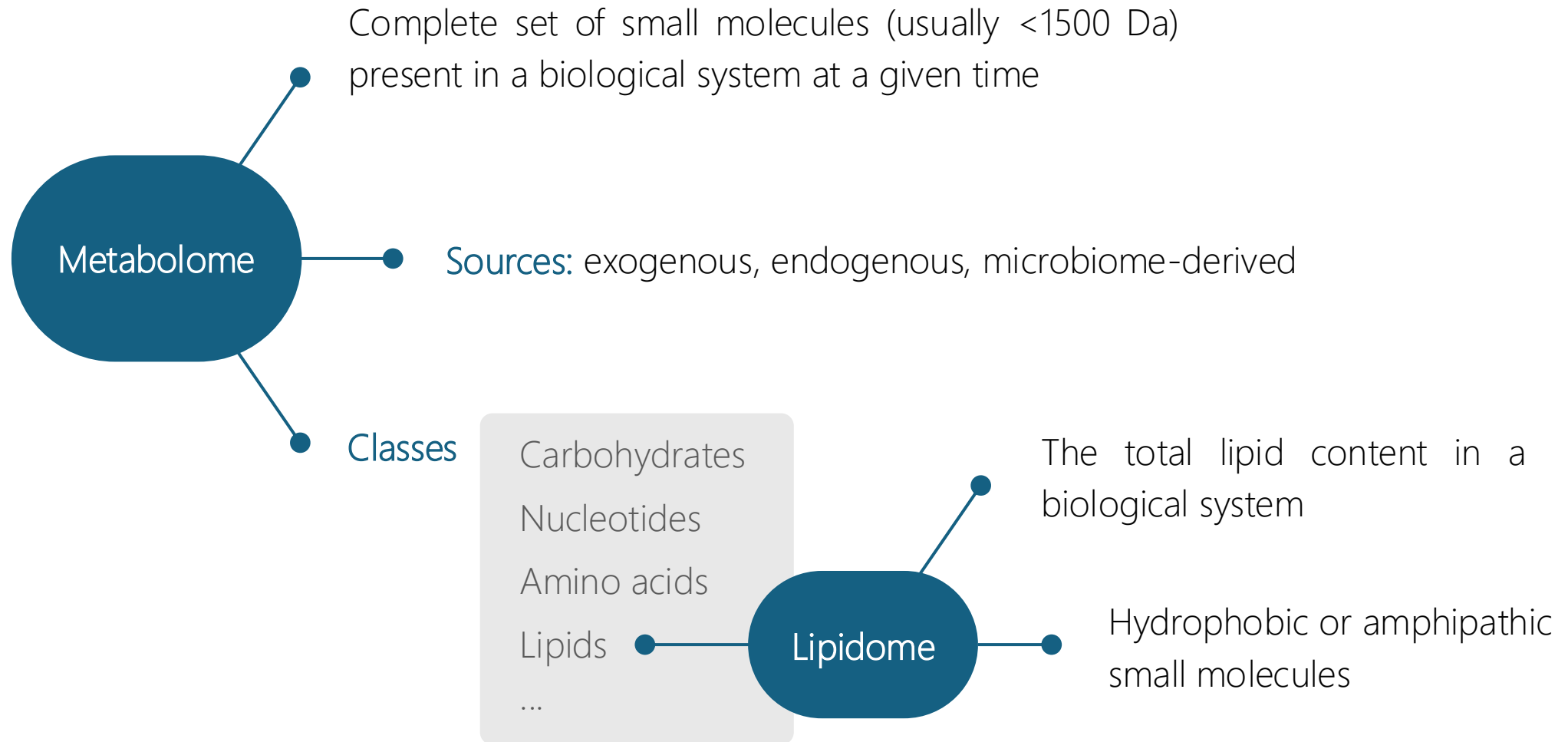
# Introduction to metabolomics



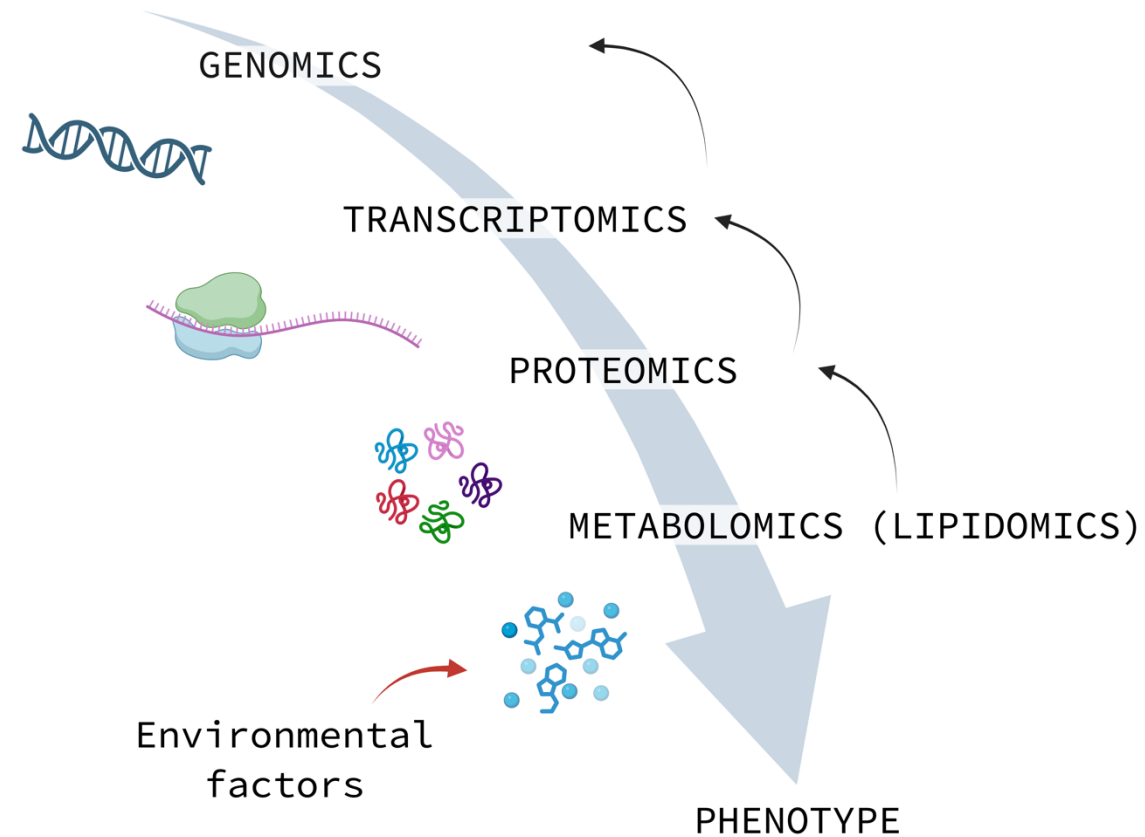
# Introduction to metabolomics



# Introduction to metabolomics



# Introduction to metabolomics



# Introduction to metabolomics

## Applications

- 1 Discovery of new disease biomarkers
- 2 The understanding of molecular mechanisms in pathophysiological processes
- 3 Nutrition research
- 4 Drug discovery
- 5 Fermentation
- 6 Environmental Monitoring and Ecotoxicology

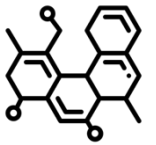
# Introduction to metabolomics

## Applications examples

Proposed biosynthetic pathway for tetrangulol in *Streptomyces* sp. KL110A



Identification of biosynthetic gene clusters (BGCs) in *Streptomyces*.



Metabolites produced by the strain



Comparison with previously characterized pathways

Trejo-Alarcon, L.M., Cano-Prieto, C., Calheiros de Carvalho, A. et al. Integrative metabolo-genomics suggests a biosynthetic pathway for tetrangulol in *Streptomyces* sp. KL110A. *World J Microbiol Biotechnol* (2025).

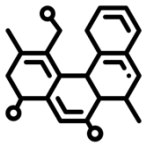
# Introduction to metabolomics

## Applications examples

Proposed biosynthetic pathway for tetrangulol in *Streptomyces* sp. KL110A



Identification of biosynthetic gene clusters (BGCs) in *Streptomyces*.



Metabolites produced by the strain



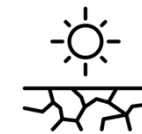
Comparison with previously characterized pathways

How do different wheat-related species respond metabolically to drought?

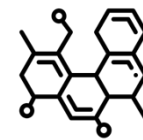
7 *Triticeae* species (drought tolerant types and cultivated types) were studied.



Normal conditions



Drought stress conditions



Sugars, aa and organic acids were consistently increased (stress protection/osmotic adjustment).

Trejo-Alarcon, L.M., Cano-Prieto, C., Calheiros de Carvalho, A. et al. Integrative metabolo-genomics suggests a biosynthetic pathway for tetrangulol in *Streptomyces* sp. KL110A. *World J Microbiol Biotechnol* (2025).

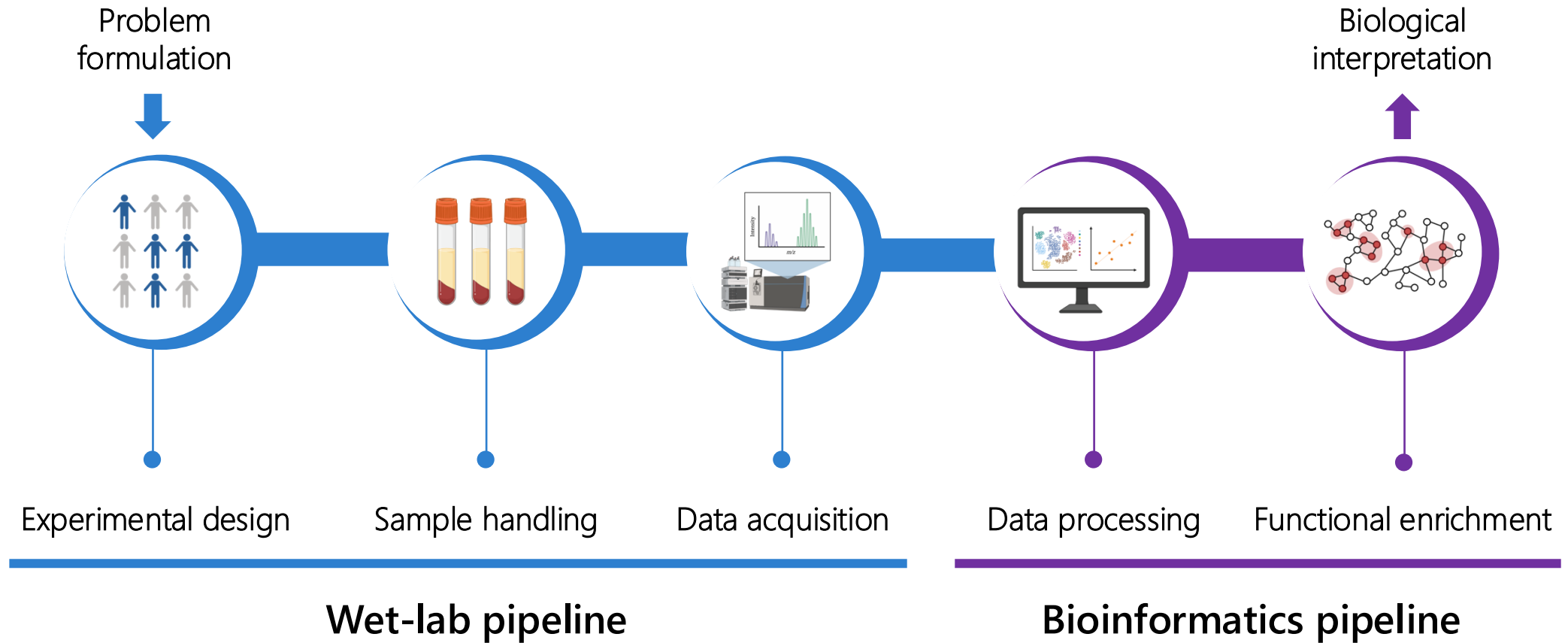
Ullah N, Yüce M, Neslihan Öztürk Gökçe Z, Budak H. Comparative metabolite profiling of drought stress in roots and leaves of seven *Triticeae* species. *BMC Genomics* (2017).

# Introduction to metabolomics

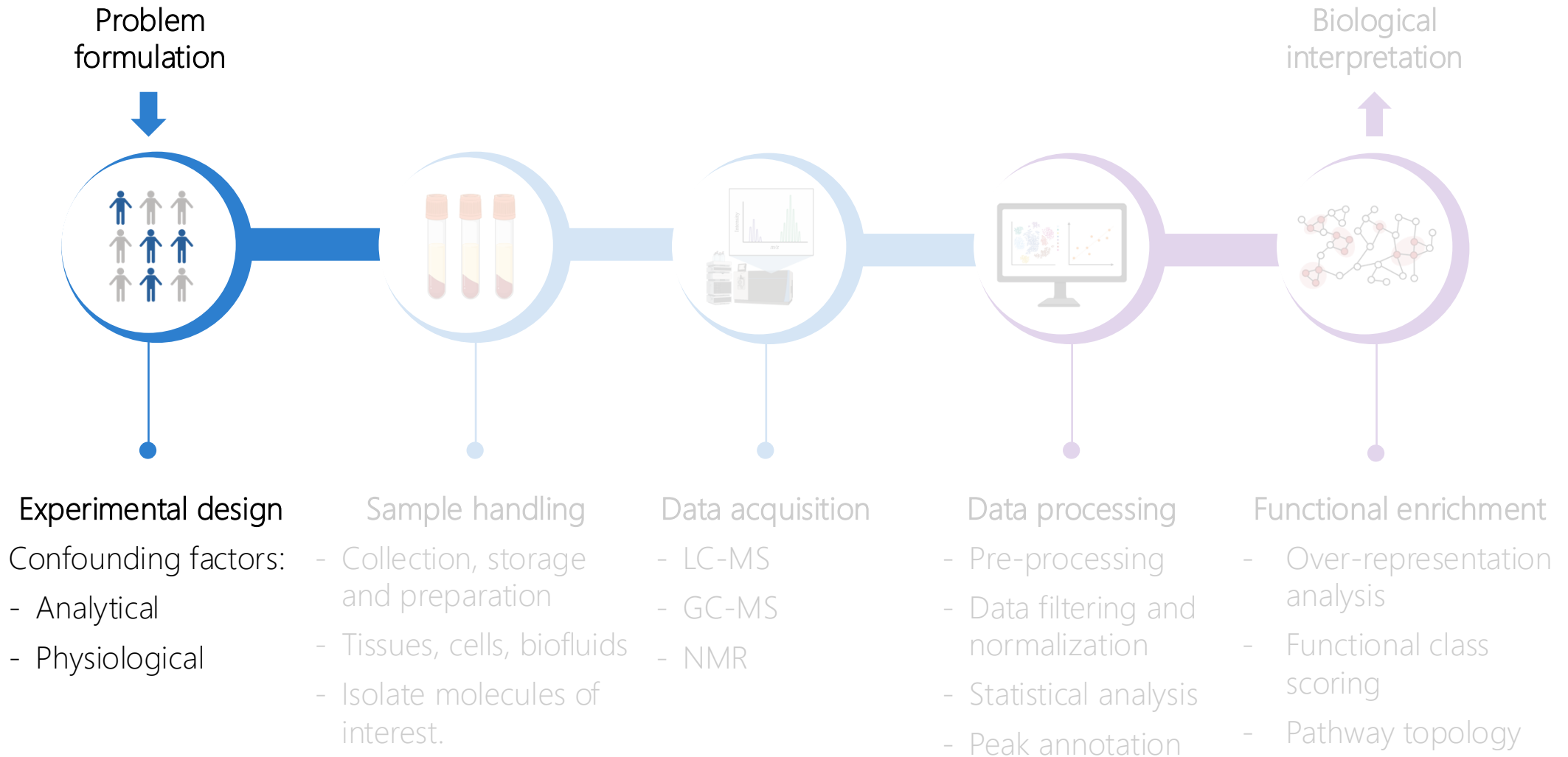
## Study types

- 1 **Untargeted** – the goal is to characterise the broadest range of metabolites present in a sample. It allows to discover unanticipated alterations; however, it results in complex datasets.
- 2 **Semi-targeted** – the goal is to identify and quantify a large set of small molecules. It allows the characterization of a pathway or some specific classes of small molecules (e.g., aminoacids).
- 3 **Targeted** – the goal is to determine absolute concentrations of a predefined set of metabolites (selected based on prior knowledge). It is usually used to validate untargeted or semi-targeted metabolomics studies.

# Metabolomics pipeline



# Metabolomics pipeline



# Metabolomics pipeline

## Experimental design

Deal with confounding factors that could correlate or mask the biology under-study:

- Analytical: batch effects, injection order, etc.
- Biological: sex, age, body mass index, etc.

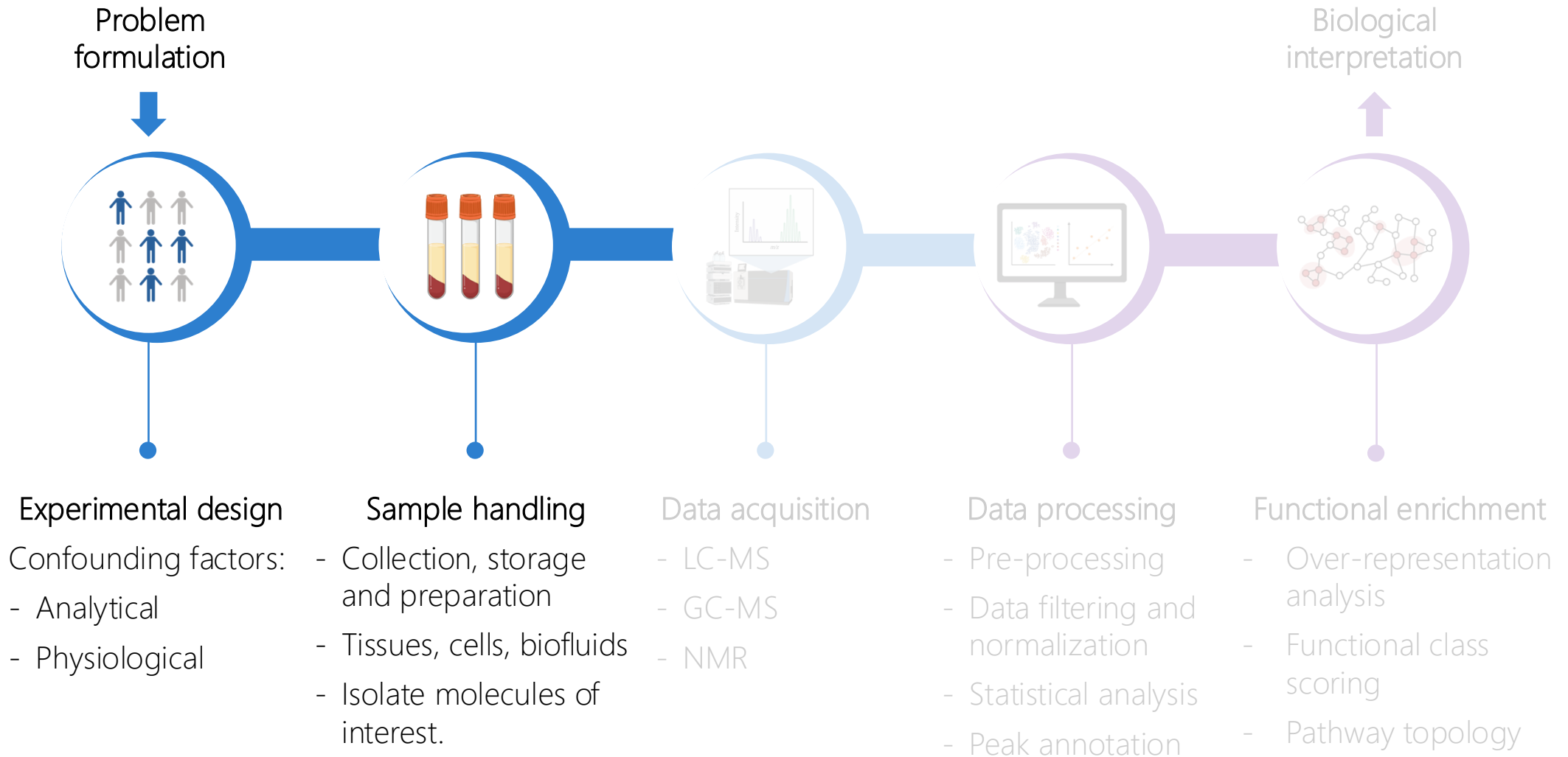


Biological samples randomization based on their characteristics

**Quality control (QC)** samples are typically prepared by taking aliquots from a pool of all the study samples.

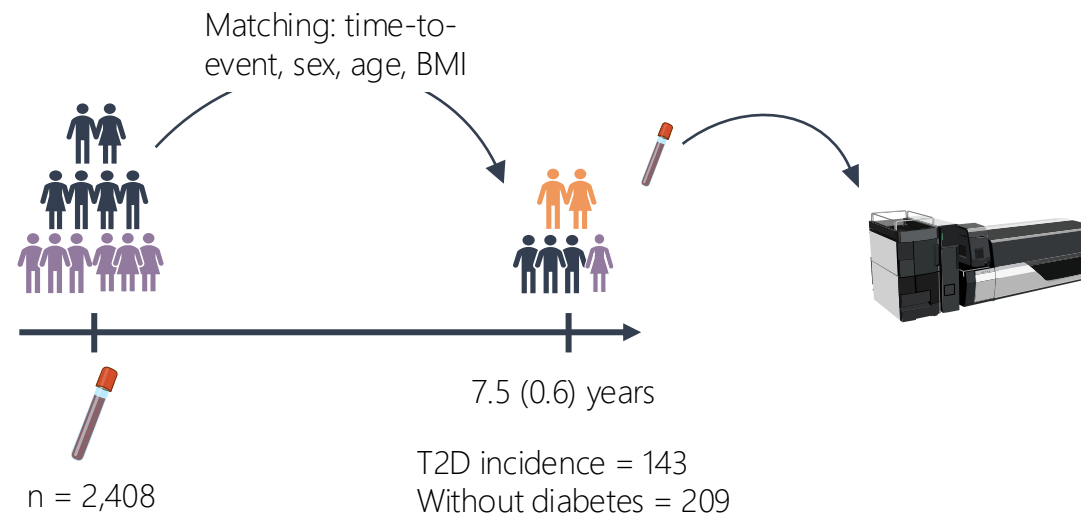
QC samples are injected at the start and periodically throughout the acquisition to stabilize the signal and correct for technical variability.

# Metabolomics pipeline



# Metabolomics pipeline

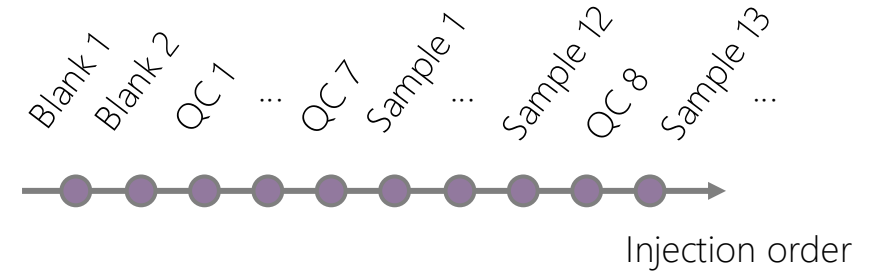
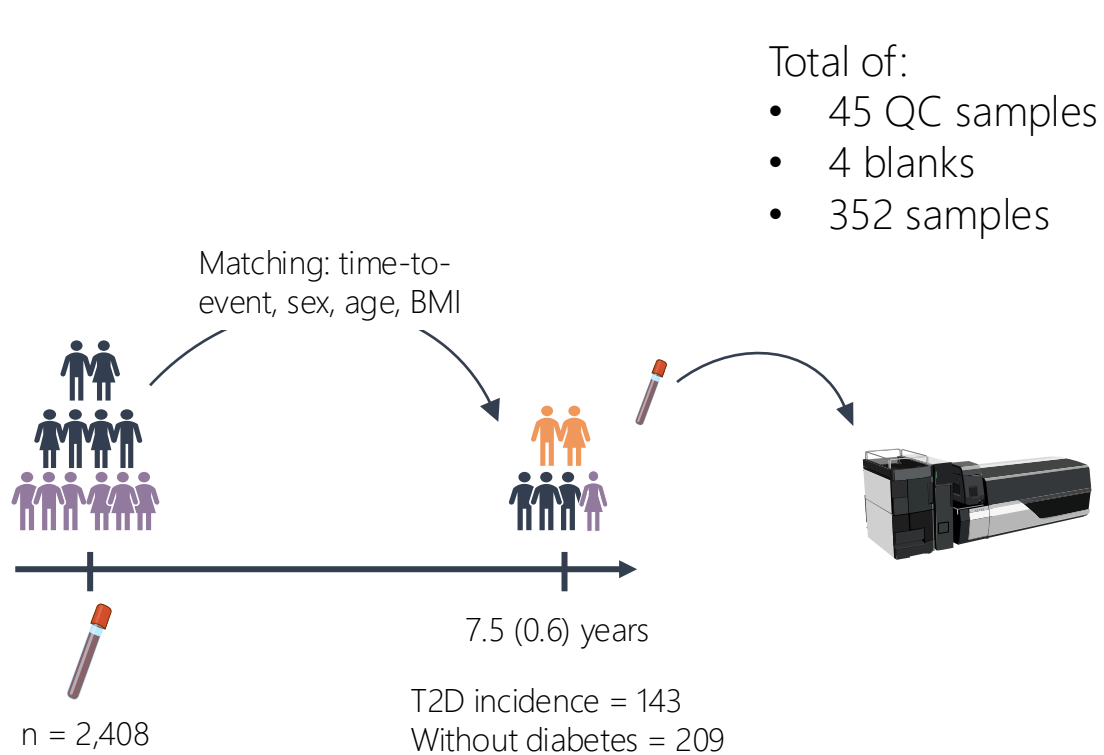
Study example – Diabetes incidence



Barranco-Altirriba et al. "Guanine and pregnenolone sulfate are associated with incident type 2 diabetes in two independent populations". *Frontiers in Endocrinology* (2025).

# Metabolomics pipeline

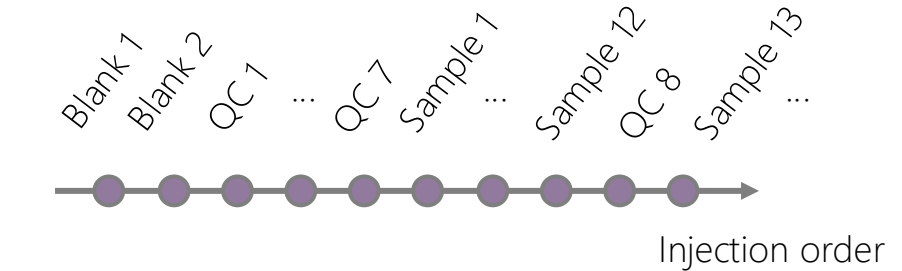
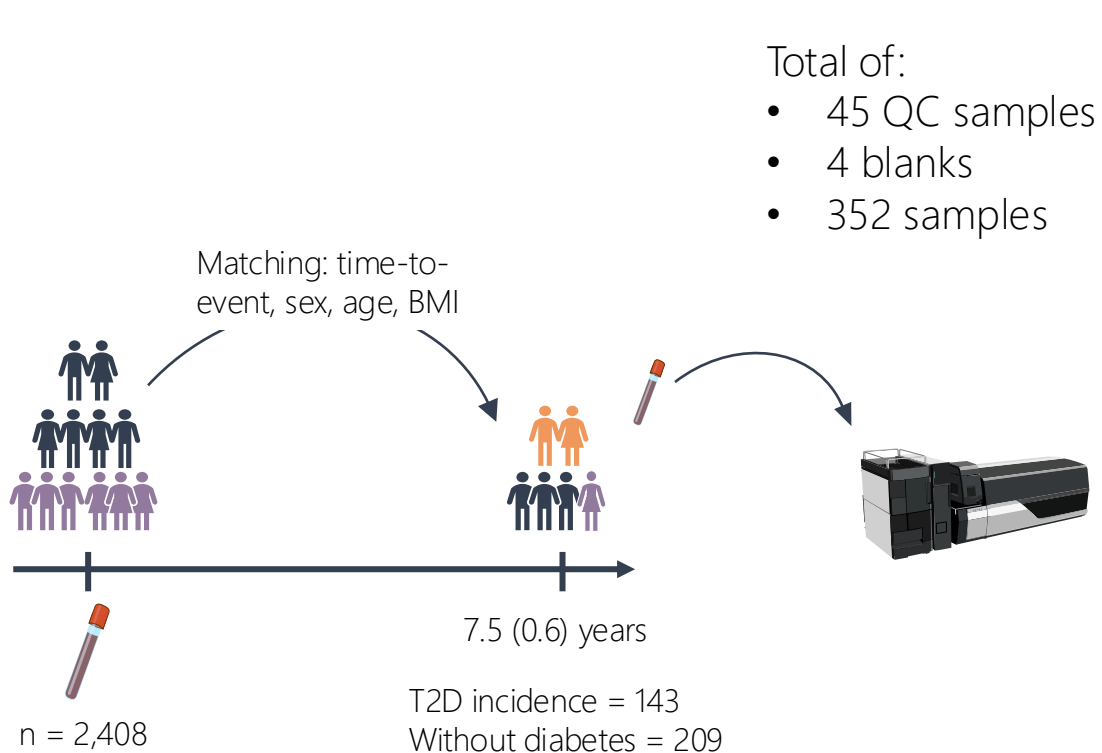
## Study example – Diabetes incidence



Barranco-Altirriba et al. "Guanine and pregnenolone sulfate are associated with incident type 2 diabetes in two independent populations". *Frontiers in Endocrinology* (2025).

# Metabolomics pipeline

## Study example – Diabetes incidence

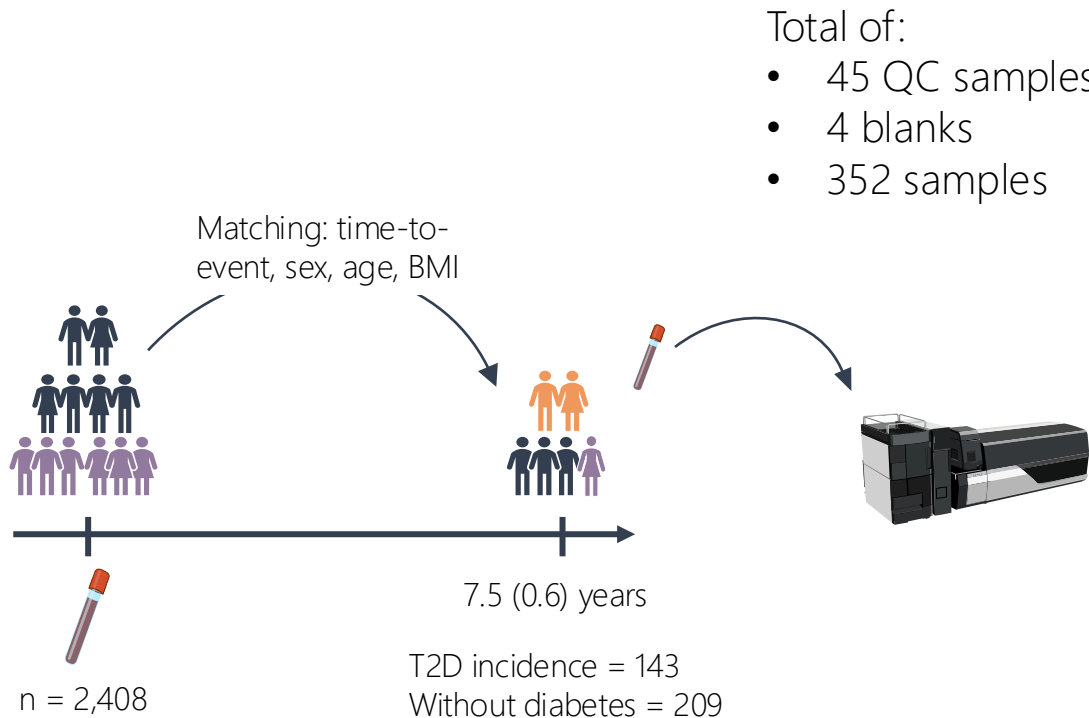


	Without T2D (N=286)	With T2D (N=143)	P-value
Sex (women)	157 (54.9%)	79 (55.2%)	
Age (years)	57.5 (12.8)	56.6 (12.0)	0.442
BMI (kg/m <sup>2</sup> )	30.5 (4.56)	30.8 (4.91)	0.483
HT (yes)	160 (55.9%)	85 (59.4%)	0.558
Insulin Resistance	2.38 (1.49)	3.05 (1.83)	<0.001
Glucose (mg/dL)	95.0 (10.6)	103 (13.2)	<0.001
Family history (yes)	134 (46.9%)	90 (63.4%)	0.002

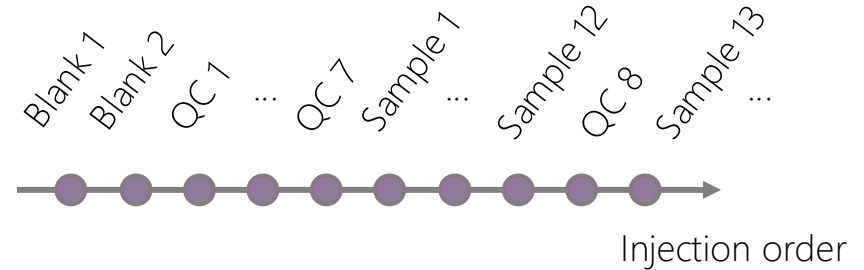
Barranco-Altirriba et al. "Guanine and pregnenolone sulfate are associated with incident type 2 diabetes in two independent populations". *Frontiers in Endocrinology* (2025).

# Metabolomics pipeline

## Study example – Diabetes incidence

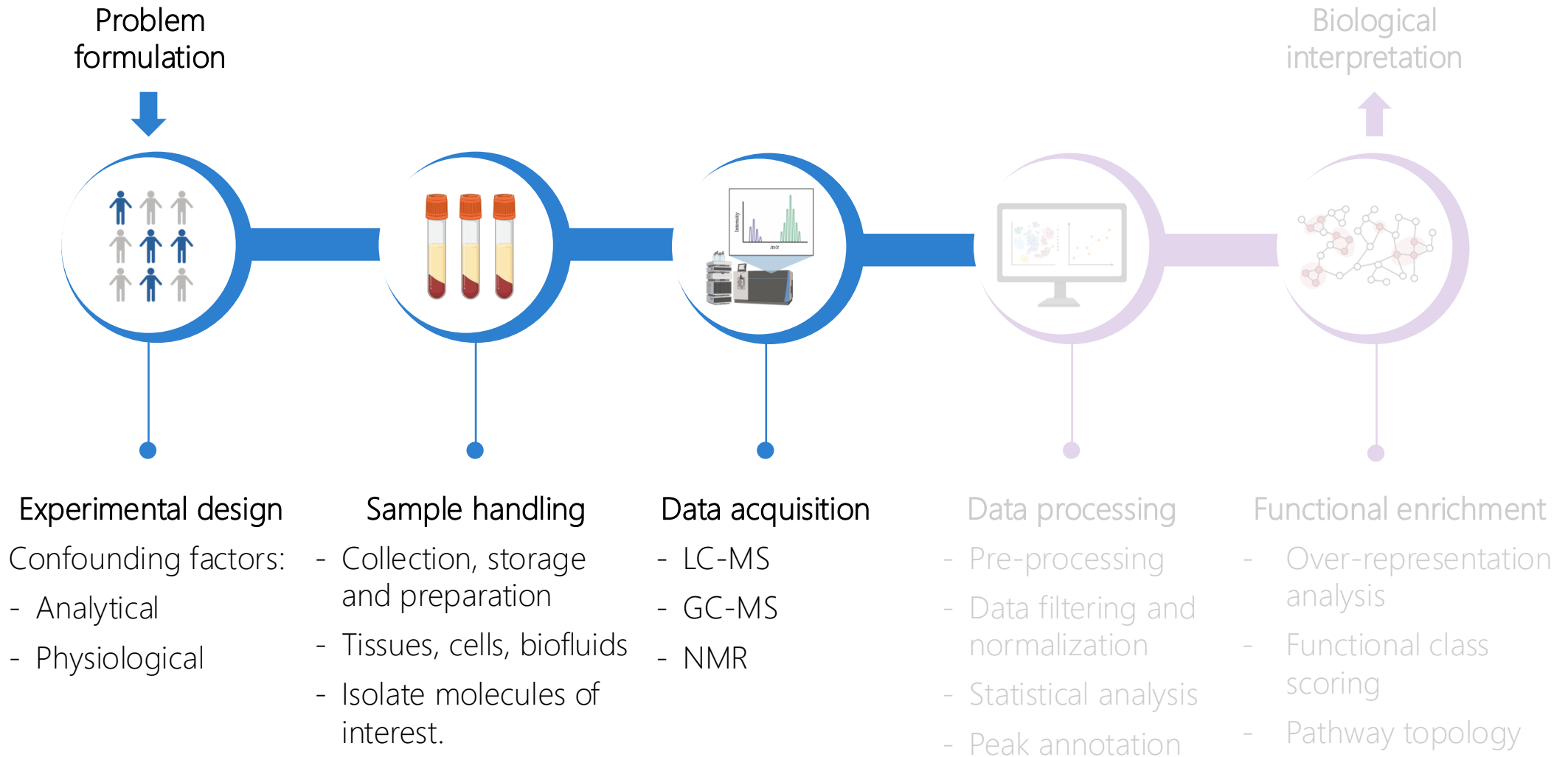


- Total of:
- 45 QC samples
  - 4 blanks
  - 352 samples



	Without T2D (N=286)	With T2D (N=143)	P-value
Sex (women)	157 (54.9%)	79 (55.2%)	
Age (years)	57.5 (12.8)	56.6 (12.0)	0.442
BMI (kg/m <sup>2</sup> )	30.5 (4.56)	30.8 (4.91)	0.483
HT (yes)	160 (55.9%)	85 (59.4%)	0.558
Insulin Resistance	2.38 (1.49)	3.05 (1.83)	<0.001
Glucose (mg/dL)	95.0 (10.6)	103 (13.2)	<0.001
Family history (yes)	134 (46.9%)	90 (63.4%)	0.002

# Metabolomics pipeline



# Metabolomics data acquisition

- Nuclear Magnetic Resonance (NMR)
- Gas chromatography-Mass spectrometry (GC-MS)
- Liquid chromatography-Mass spectrometry (LC-MS)

# Metabolomics data acquisition

Nuclear Magnetic Resonance (NMR)

## STRENGTHS

- 1 Non-destructive
- 2 Easily quantifiable
- 3 No prior separation required
- 4 Identification of novel compounds
- 5 High reproducibility

## LIMITATIONS

- 1 Low-sensitivity
- 2 High detection limit
- 3 Large sample volumes

# Metabolomics data acquisition

Gas chromatography-Mass spectrometry (GC-MS)

## STRENGTHS

- 1 More mature and cost-effective
- 2 Excellent separation reproducibility
- 3 Availability of spectral libraries facilitates metabolite identification

## LIMITATIONS

- 1 Sample derivatization required
- 2 Smaller coverage
- 3 Longer experimental times

# Metabolomics data acquisition

Liquid chromatography-Mass spectrometry (LC-MS)

## STRENGTHS

- 1 Broader range of metabolites
- 2 Simpler sample preparation
- 3 High sensitivity

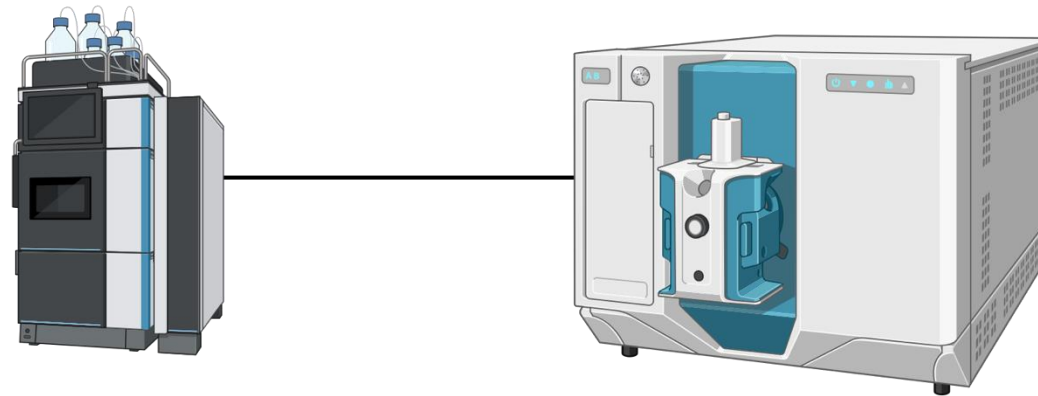
## LIMITATIONS

- 1 Higher cost
- 2 Lower reproducibility
- 3 Reduced compatibility with volatile compounds

**Preferred** choice in metabolomics: **more than 70%** of the metabolomics studies were using LC-MS in 2022 .

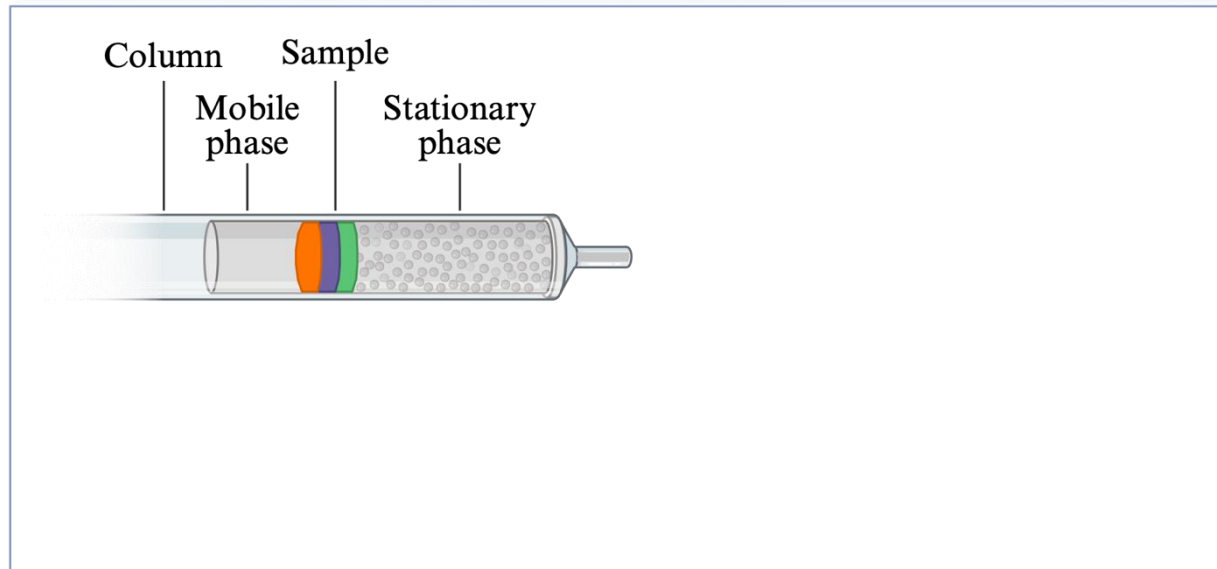
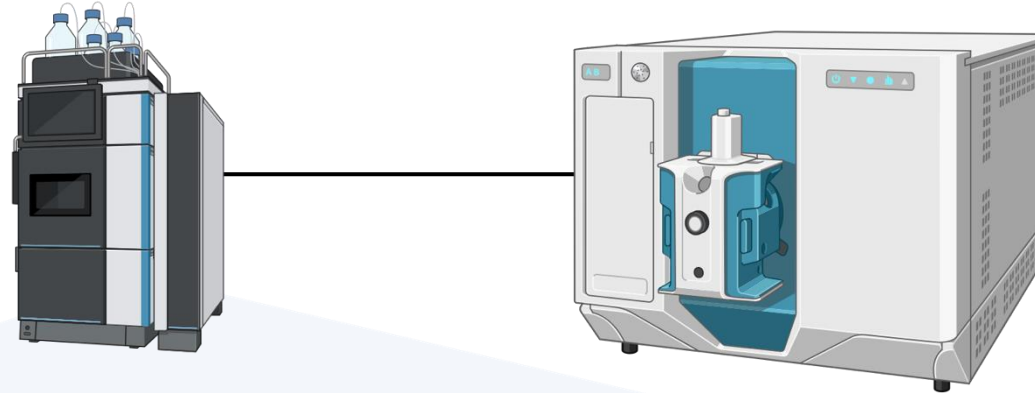
# Metabolomics data acquisition

Liquid chromatography-Mass spectrometry (LC-MS/MS)



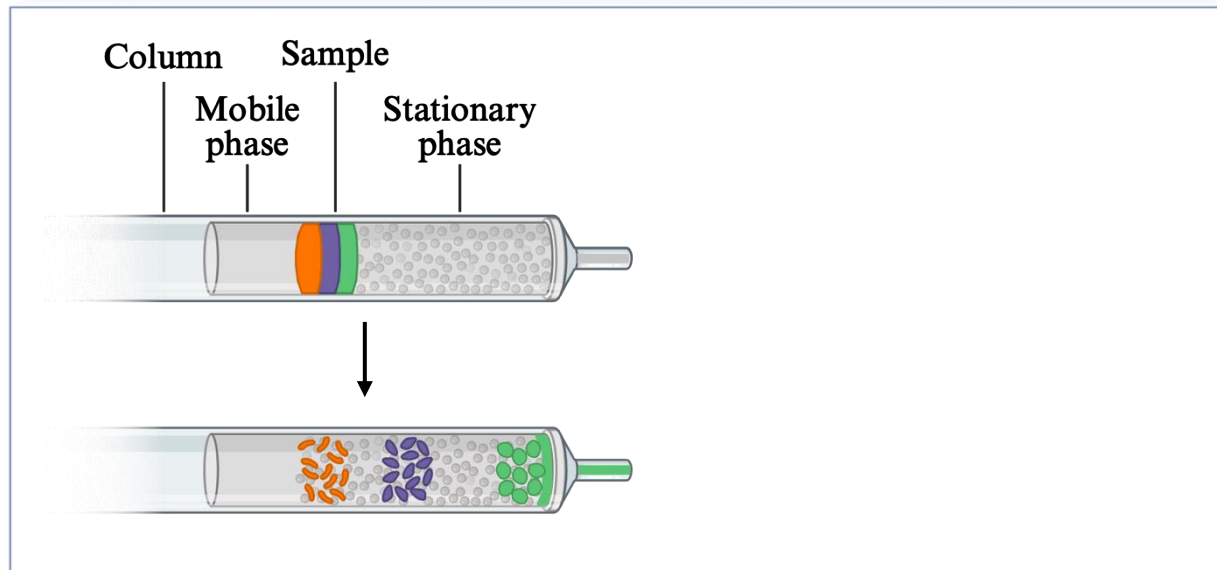
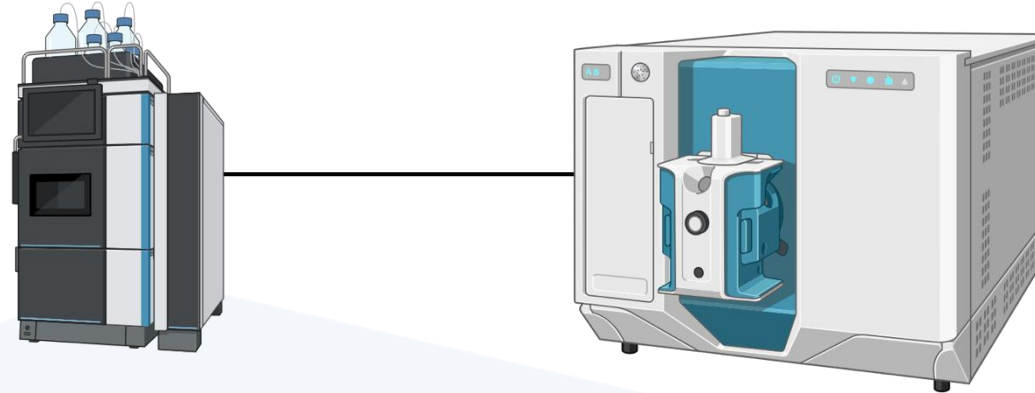
# Metabolomics data acquisition

Liquid chromatography-Mass spectrometry (LC-MS/MS)



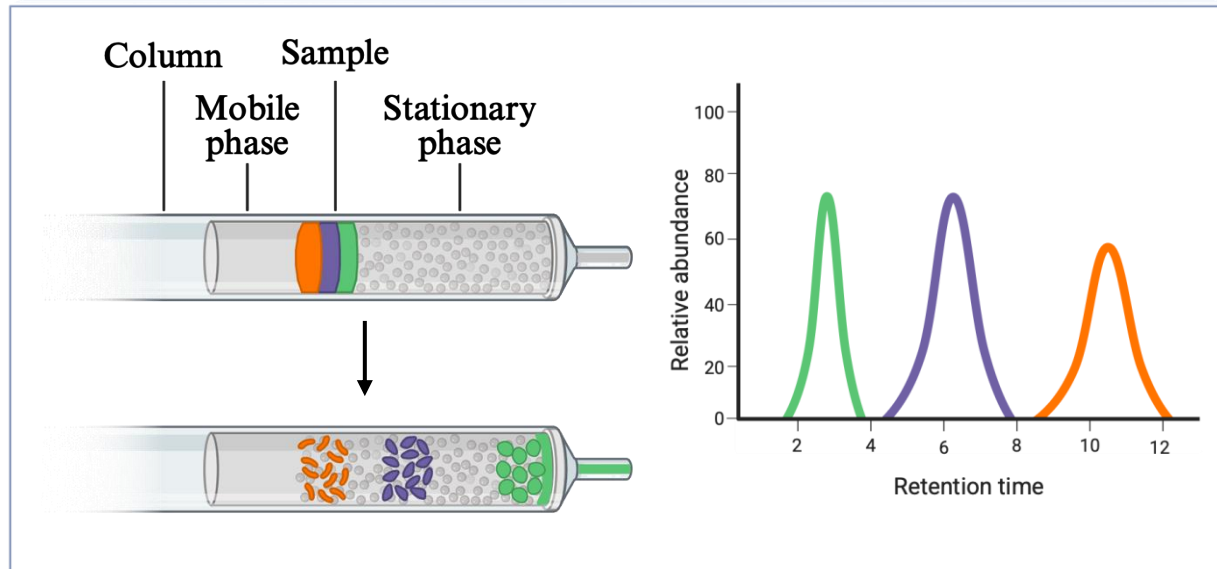
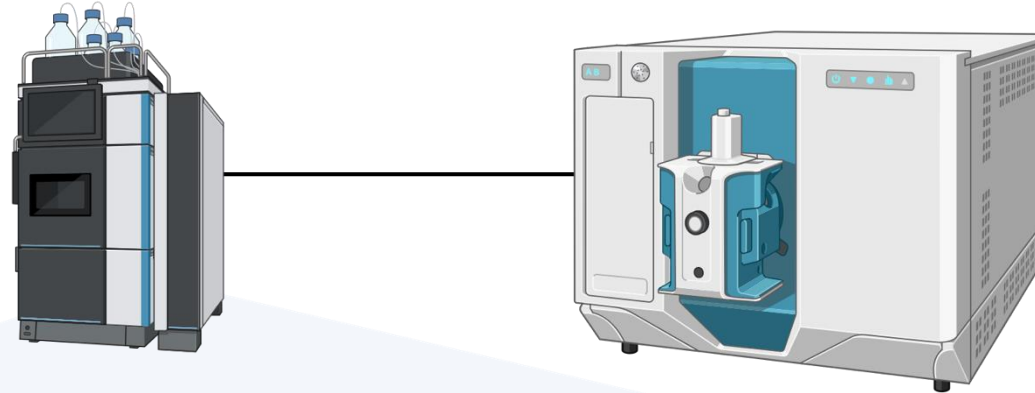
# Metabolomics data acquisition

Liquid chromatography-Mass spectrometry (LC-MS/MS)



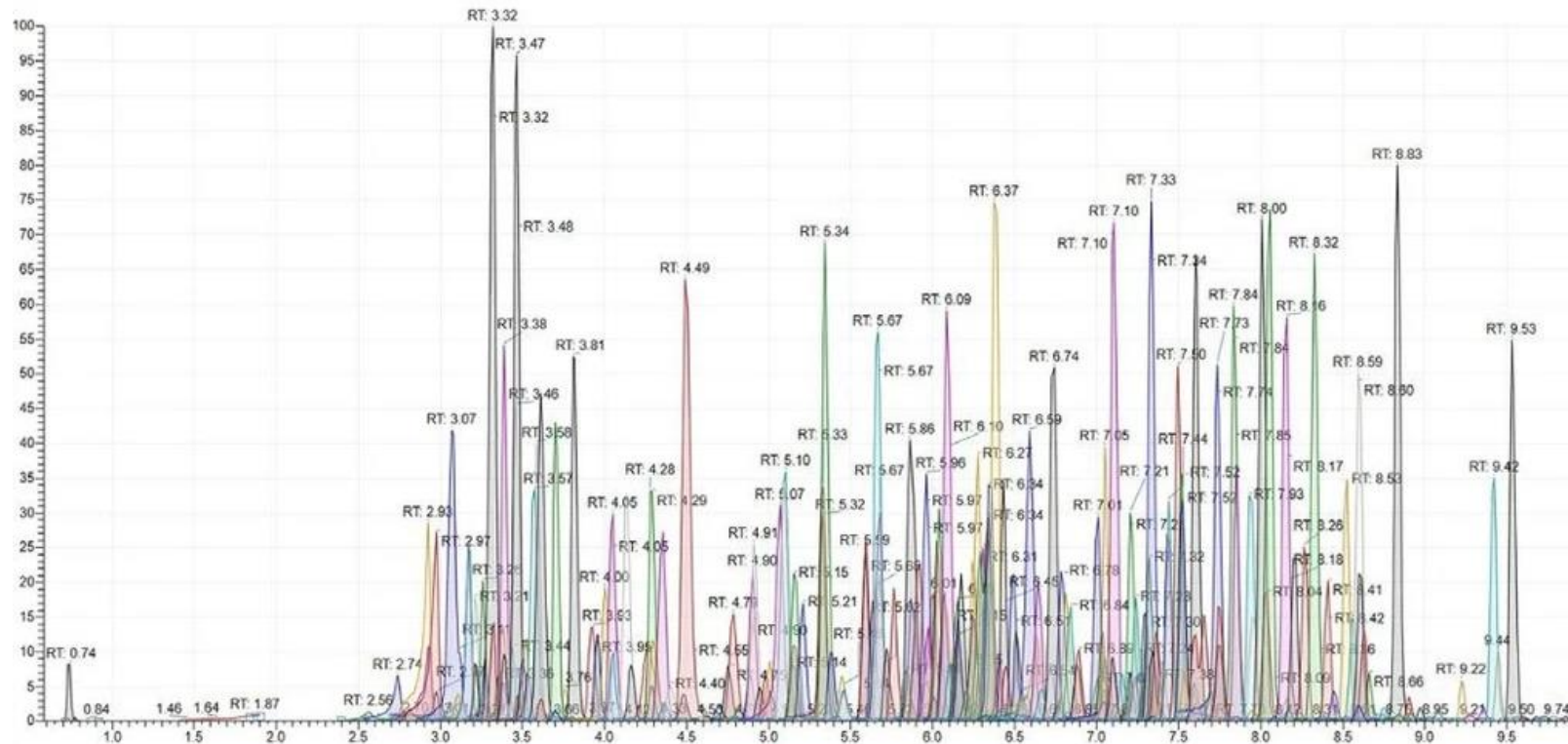
# Metabolomics data acquisition

Liquid chromatography-Mass spectrometry (LC-MS/MS)



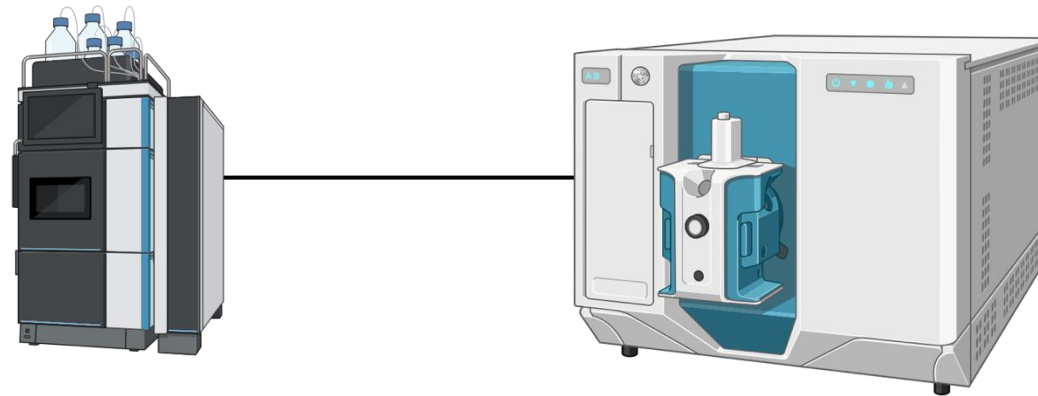
# Metabolomics data acquisition

Liquid chromatography-Mass spectrometry (LC-MS/MS)



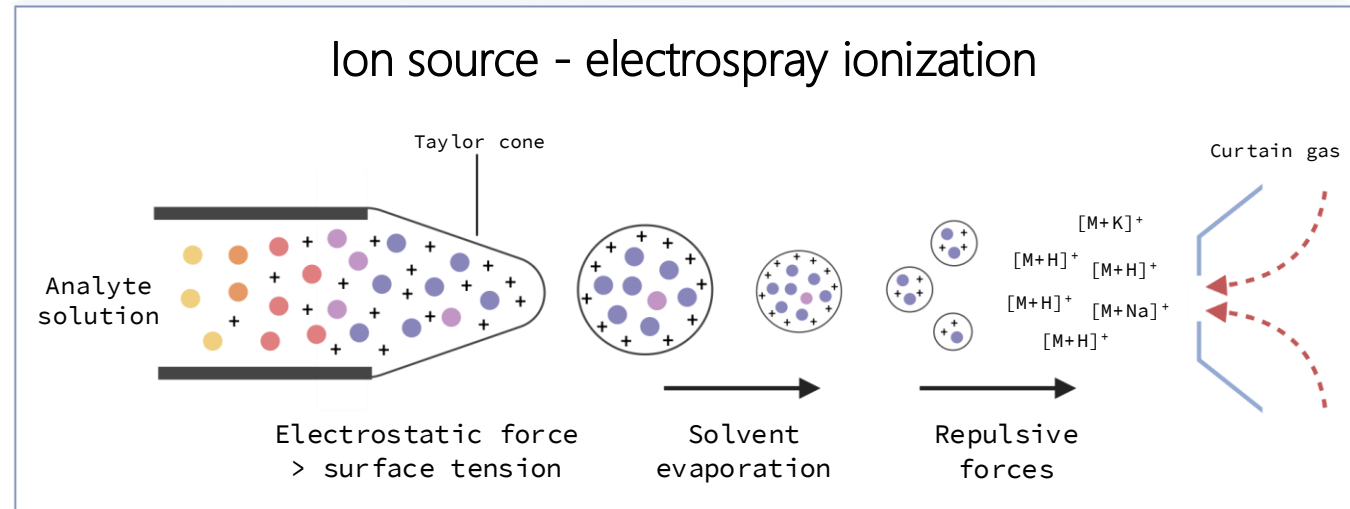
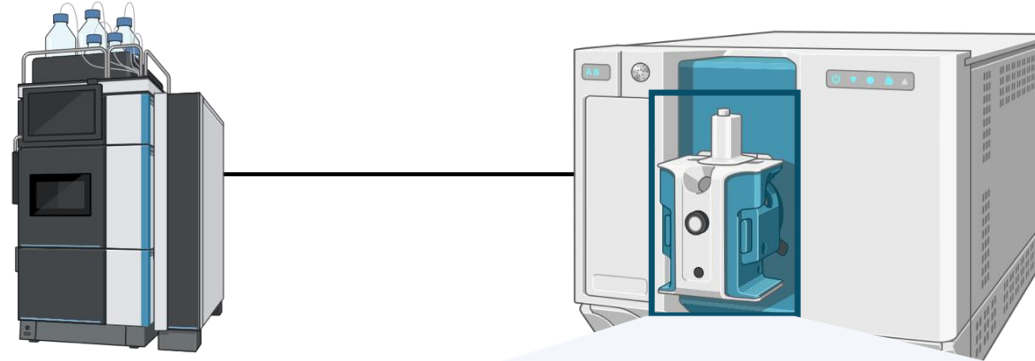
# Metabolomics data acquisition

Liquid chromatography-Mass spectrometry (LC-MS/MS)



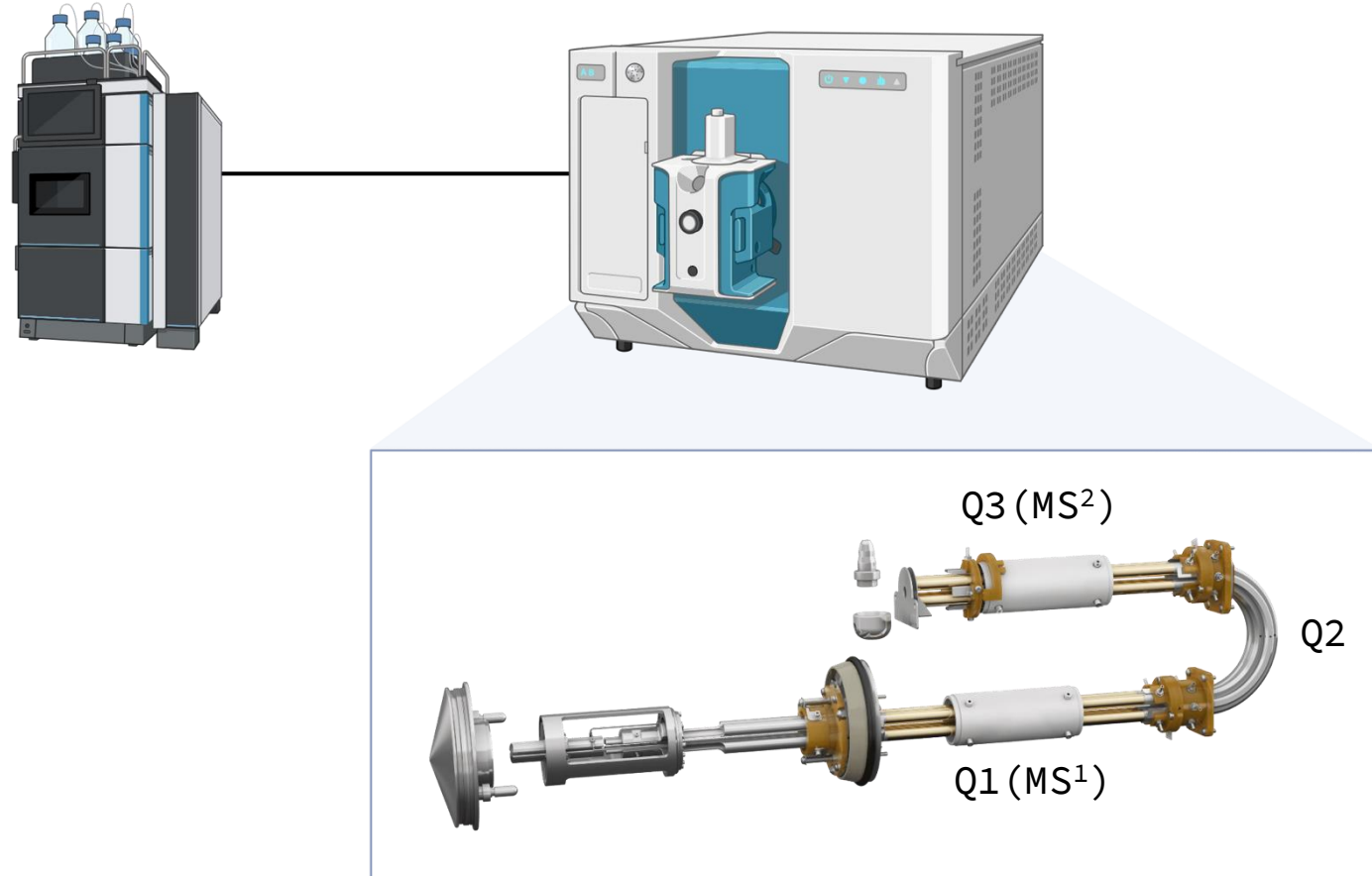
# Metabolomics data acquisition

Liquid chromatography-Mass spectrometry (LC-MS/MS)



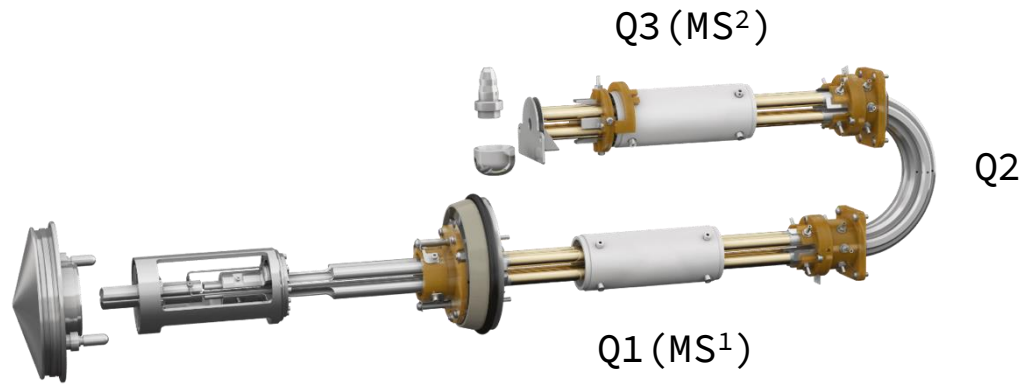
# Metabolomics data acquisition

Liquid chromatography-Mass spectrometry (LC-MS/MS)



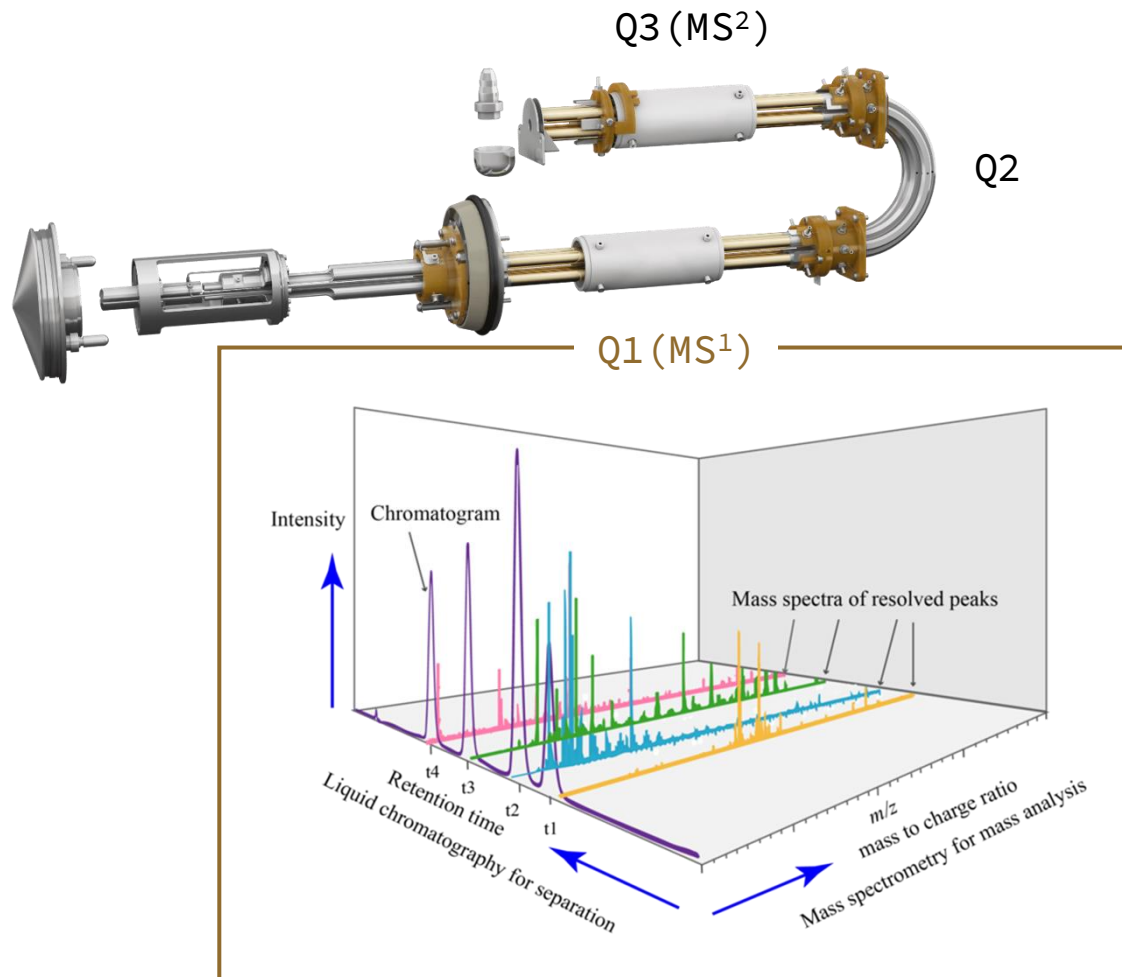
# Metabolomics data acquisition

Liquid chromatography-Mass spectrometry (LC-MS/MS)



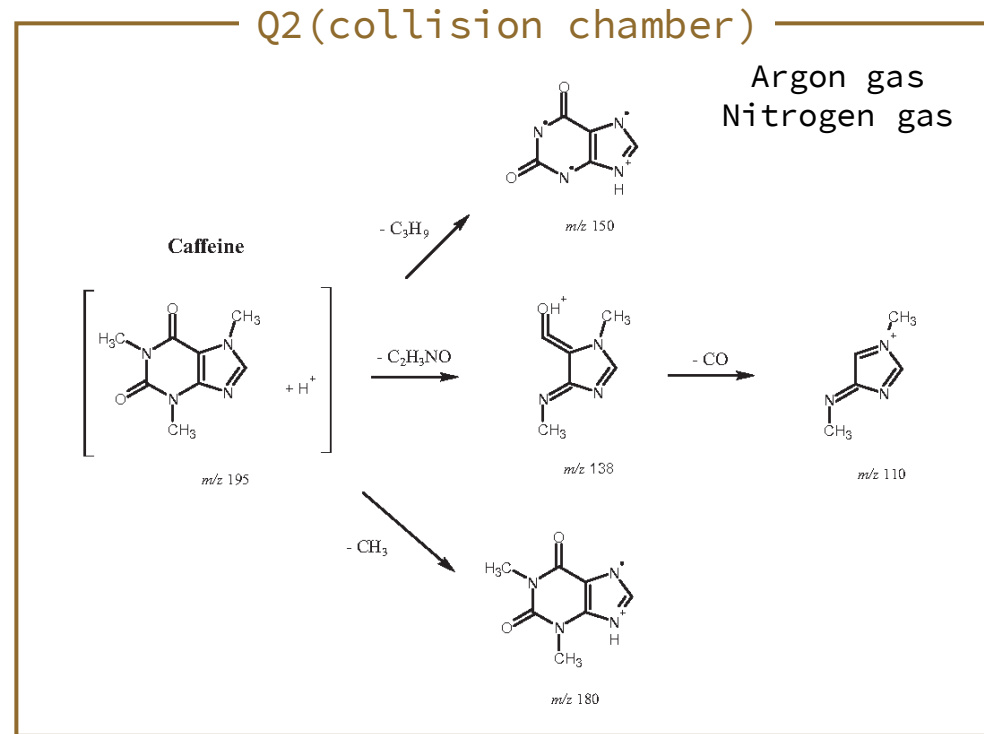
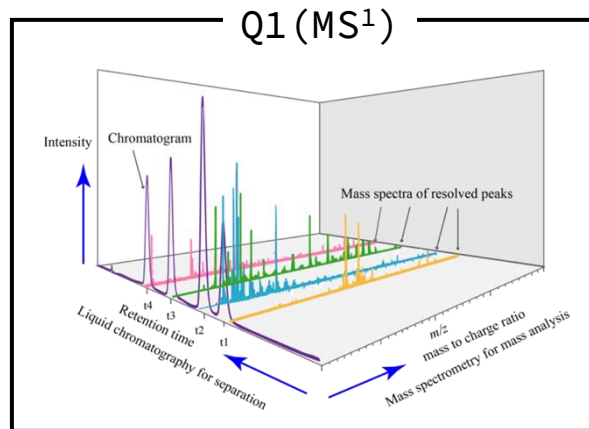
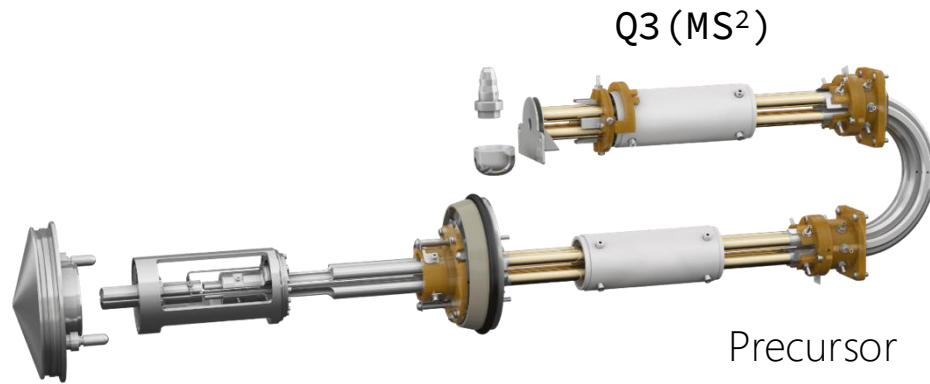
# Metabolomics data acquisition

Liquid chromatography-Mass spectrometry (LC-MS/MS)



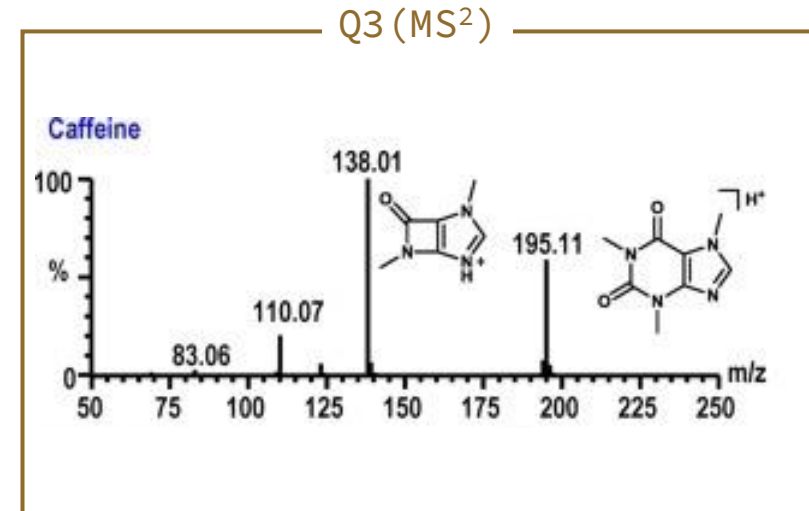
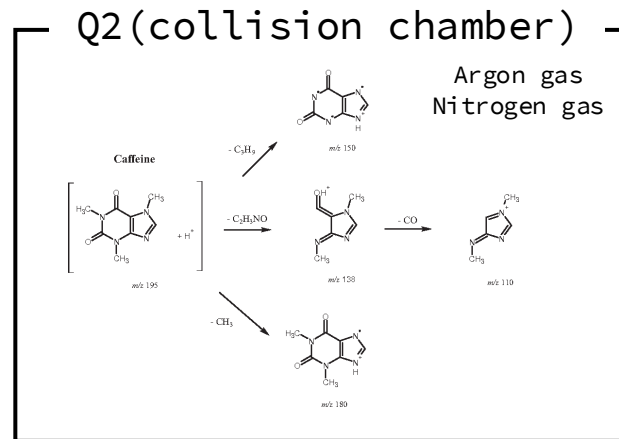
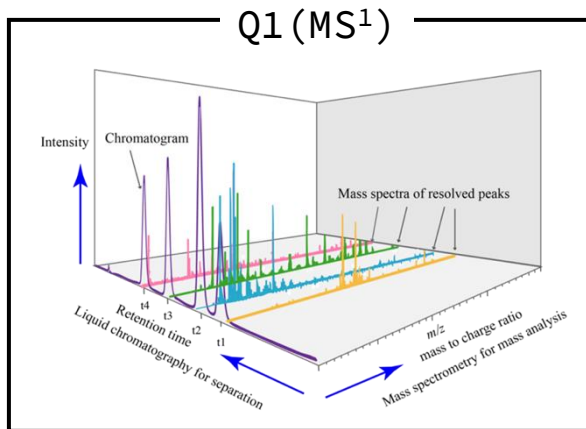
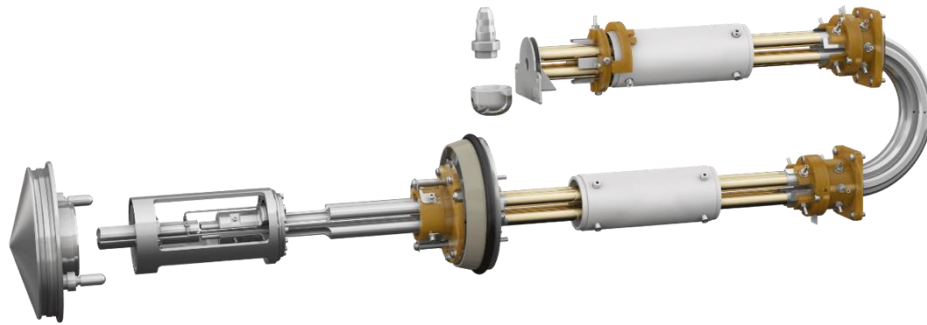
# Metabolomics data acquisition

Liquid chromatography-Mass spectrometry (LC-MS/MS)



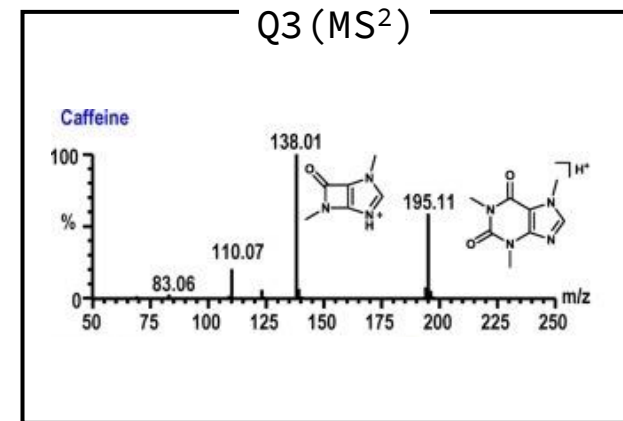
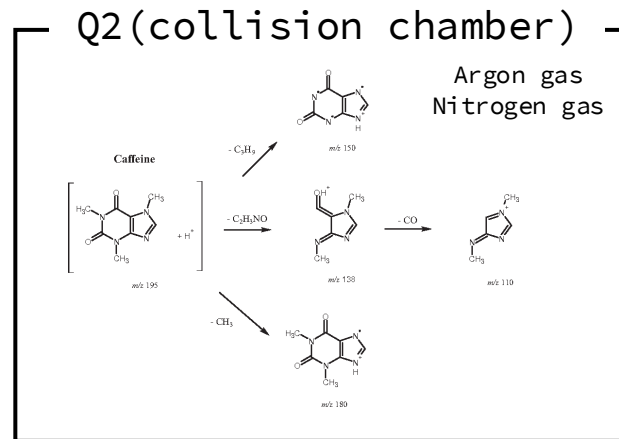
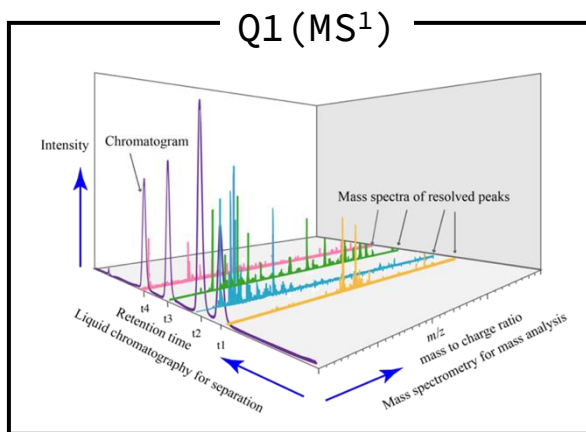
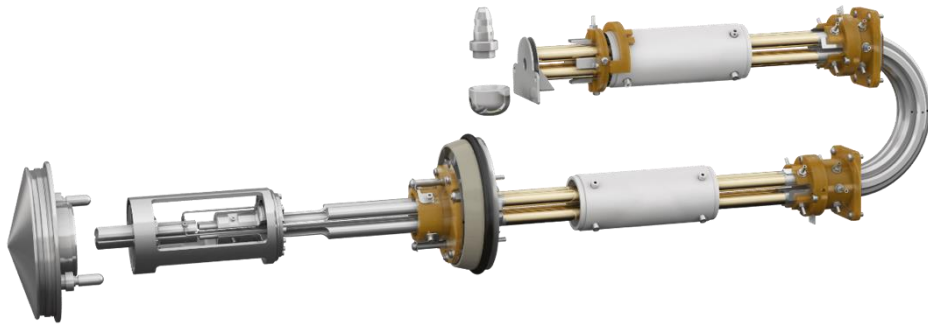
# Metabolomics data acquisition

Liquid chromatography-Mass spectrometry (LC-MS/MS)



# Metabolomics data acquisition

Liquid chromatography-Mass spectrometry (LC-MS/MS)



# Metabolomics data acquisition

Liquid chromatography-Mass spectrometry (LC-MS/MS)

## Data Dependent Acquisition (DDA)

Ions are selected based on intensity.

The MS makes real-time decisions during the run: MS<sup>1</sup> scan → Selection of Top K most intense ions

### Pros

High-quality MS/MS spectra

### Cons

Biased toward abundant metabolites

# Metabolomics data acquisition

Liquid chromatography-Mass spectrometry (LC-MS/MS)

## Data Dependent Acquisition (DDA)

Ions are selected based on intensity.

The MS makes real-time decisions during the run: MS<sup>1</sup> scan → Selection of Top K most intense ions

### Pros

High-quality MS/MS spectra

### Cons

Biased toward abundant metabolites

## Data Independent Acquisition (DIA)

It fragments all ions within predefined m/z windows.

### Pros

Much more comprehensive coverage

### Cons

MS/MS spectra are complex and overlapping

# Metabolomics data acquisition

Liquid chromatography-Mass spectrometry (LC-MS/MS)

## Data Dependent Acquisition (DDA)

Ions are selected based on intensity.

The MS makes real-time decisions during the run: MS<sup>1</sup> scan → Selection of Top K most intense ions

### Pros

High-quality MS/MS spectra

### Cons

Biased toward abundant metabolites

## Data Independent Acquisition (DIA)

It fragments all ions within predefined m/z windows.

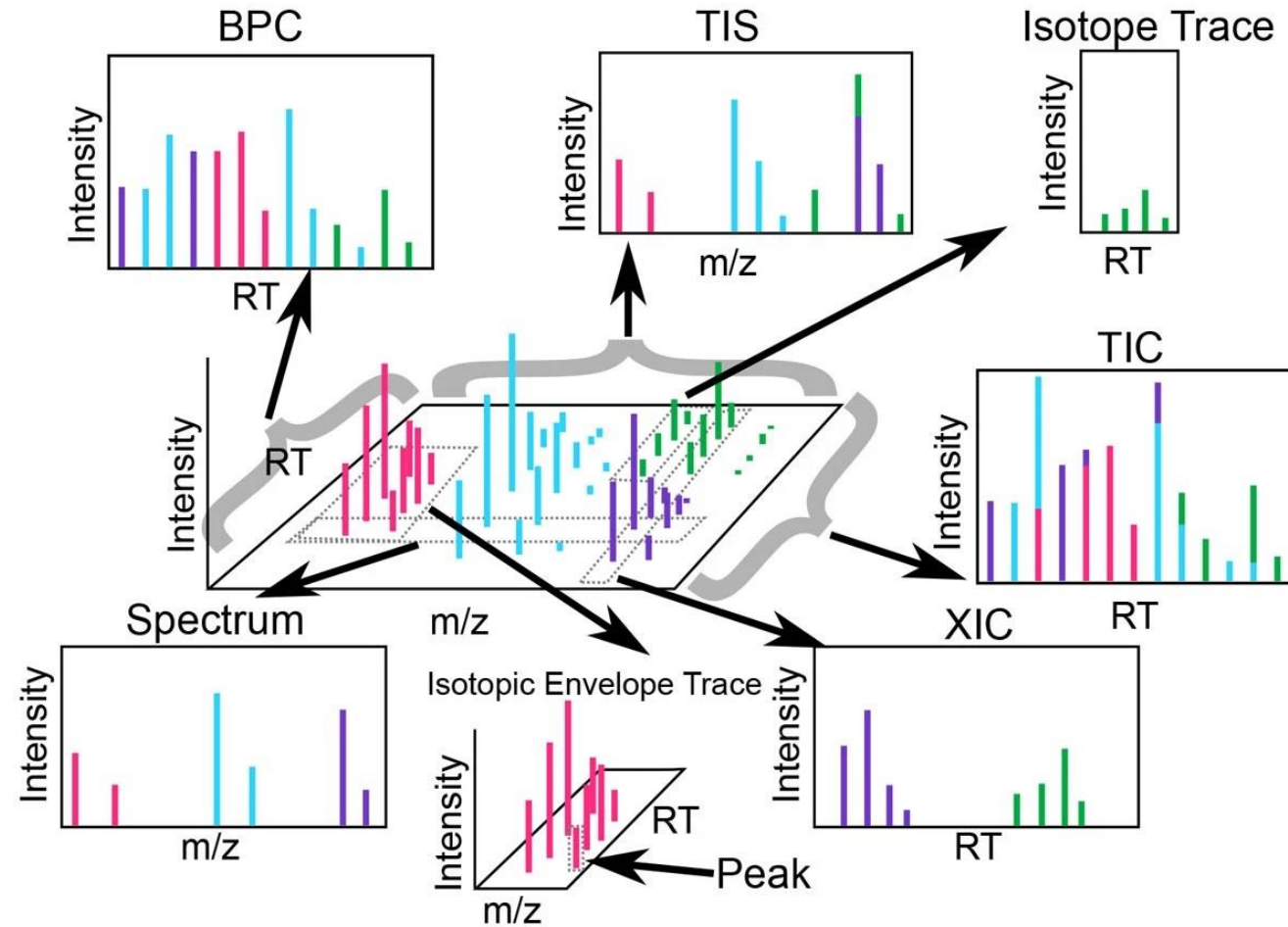
### Pros

Much more comprehensive coverage

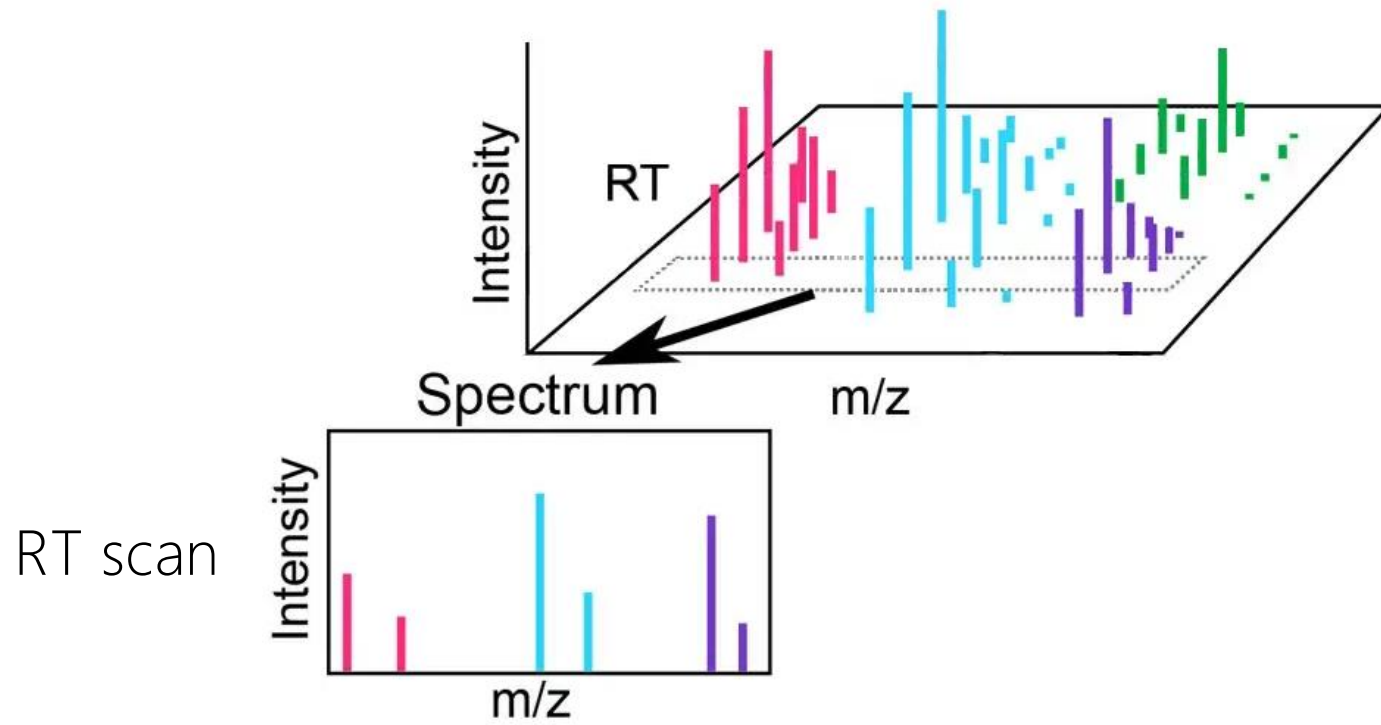
### Cons

MS/MS spectra are complex and overlapping

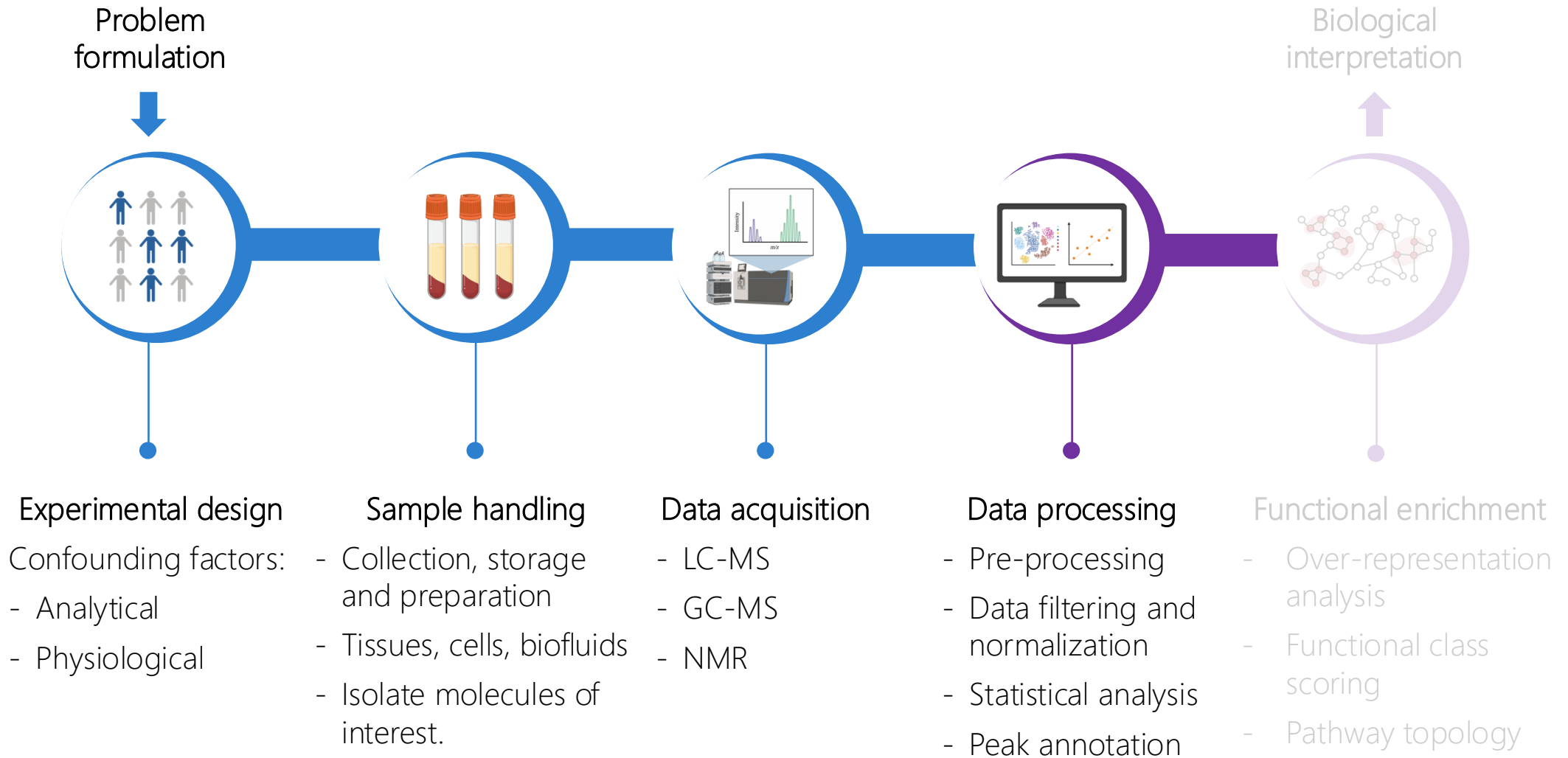
# LC-MS data



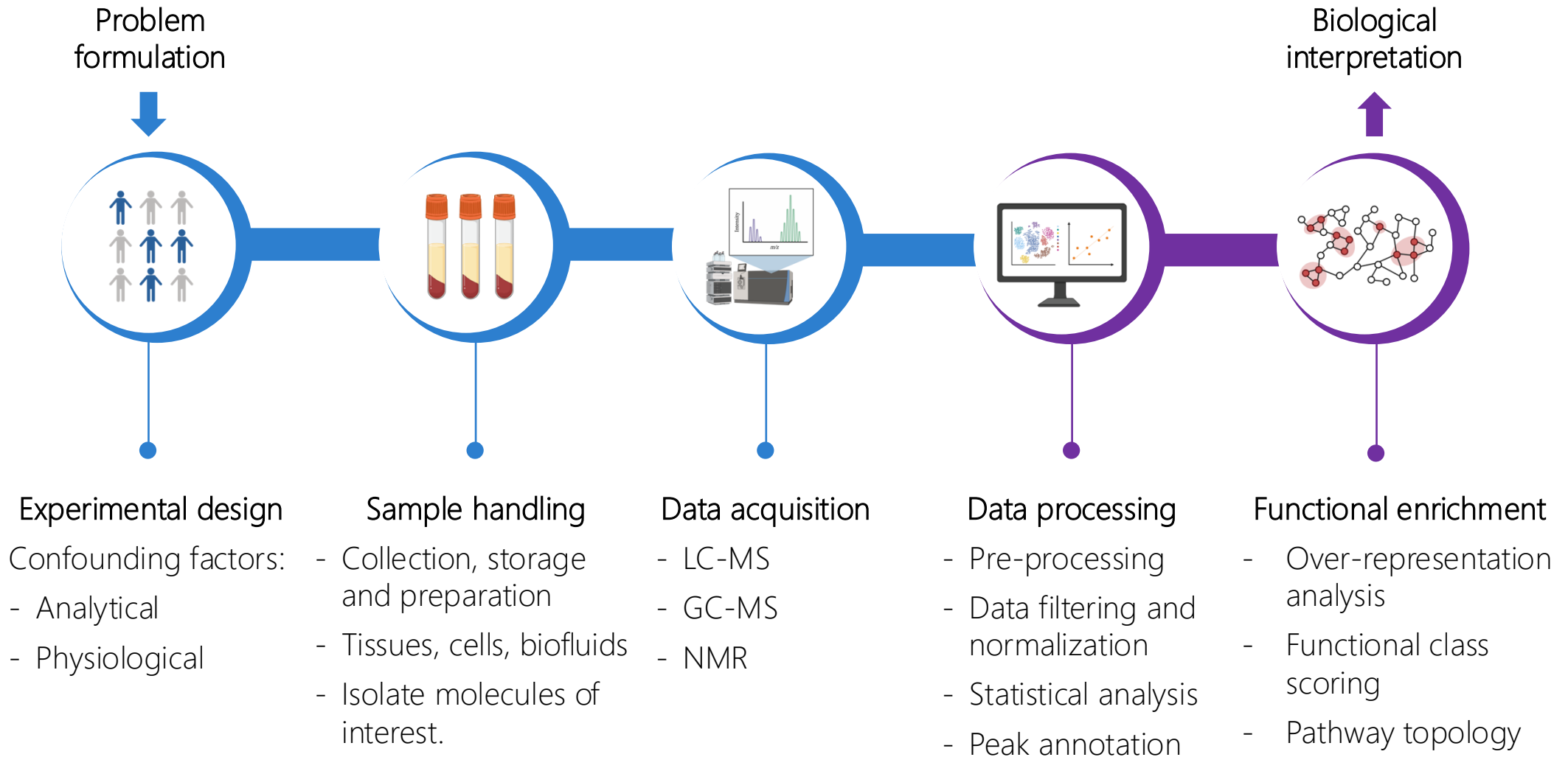
# LC-MS data



# Metabolomics pipeline

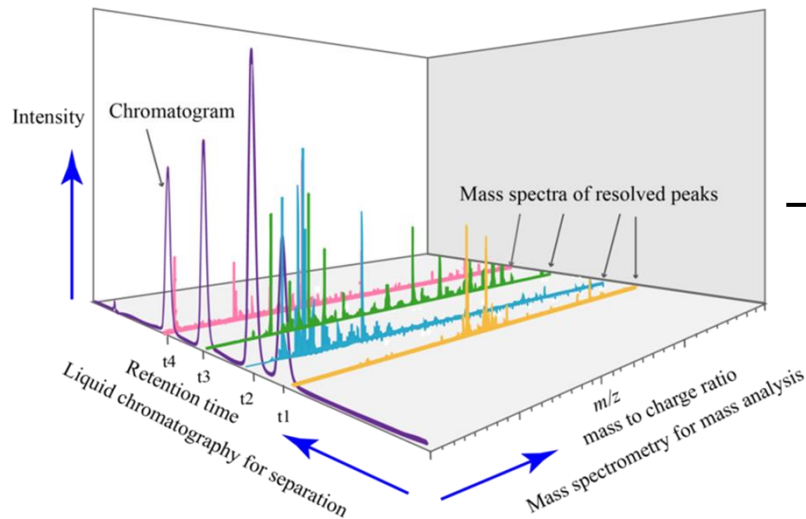


# Metabolomics pipeline



# Pre-processing

The goal is to convert **raw LC-MS data** into a **table of n peaks** defined by  $m/z$ ,  $rt$  and an intensity value for each sample  $k$ .



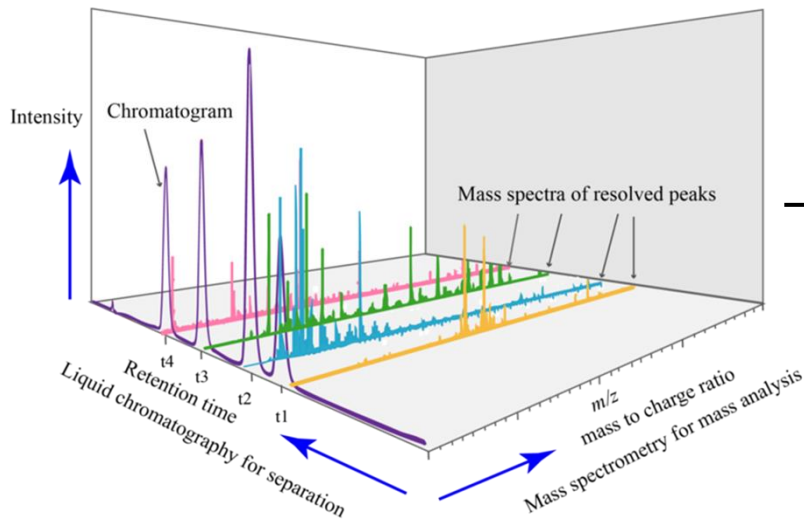
— Pre-processing →

$m/z$	RT	$I_1$	...	$I_k$
$m/z_1$	$RT_1$	$I_{11}$	...	$I_{1k}$
$m/z_2$	$RT_2$	$I_{21}$	...	$I_{2k}$
...	...	...	...	...
$m/z_n$	$RT_n$	$I_{n1}$	...	$I_{nk}$

1. Peak detection
2. Peak alignment and matching

# Pre-processing

The goal is to convert **raw LC-MS data** into a **table of  $n$  peaks** defined by  $m/z$ ,  $rt$  and an intensity value for each sample  $k$ .



— Pre-processing →

$m/z$	RT	$I_1$	...	$I_k$
$m/z_1$	$RT_1$	$I_{11}$	...	$I_{1k}$
$m/z_2$	$RT_2$	$I_{21}$	...	$I_{2k}$
...	...	...	...	...
$m/z_n$	$RT_n$	$I_{n1}$	...	$I_{nk}$

1. Peak detection
2. Peak alignment and matching



Positive acquisition mode: 2545 peaks  
Negative acquisition mode: 2187 peaks

# Pre-processing

## Centroiding

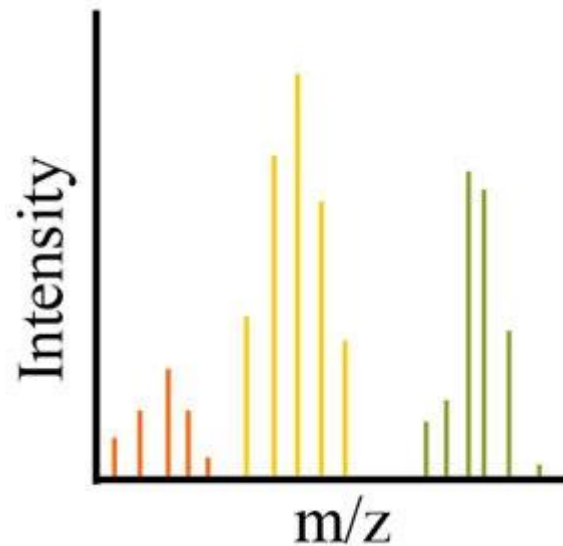
- The MS has finite resolving power → ions are detected as profile peaks around their true  $m/z$  rather than as a single exact  $m/z$
- Centroiding reduces the amount of data without much loss of information.
- The next processing steps require centroided data.
- Many manufacturers already apply centroiding of the profile data (during acquisition or immediately after).

# Pre-processing

Centroiding with OpenMS PeakPickerHiRes

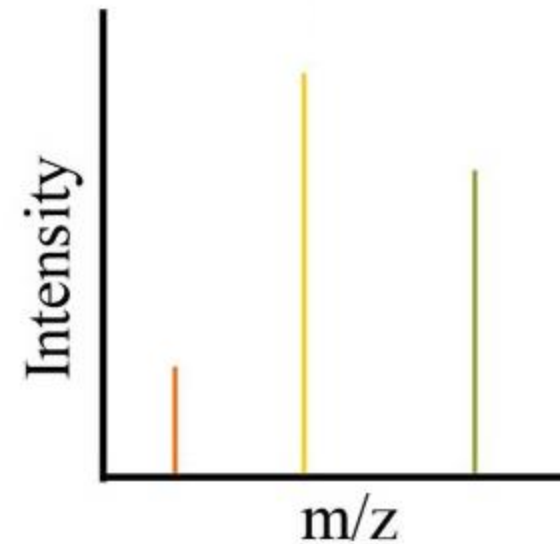
**Profile mode**

MS1 spectrum (scan at RT = 5.2 min)



**Centroid mode**

MS1 spectrum (scan at RT = 5.2 min)



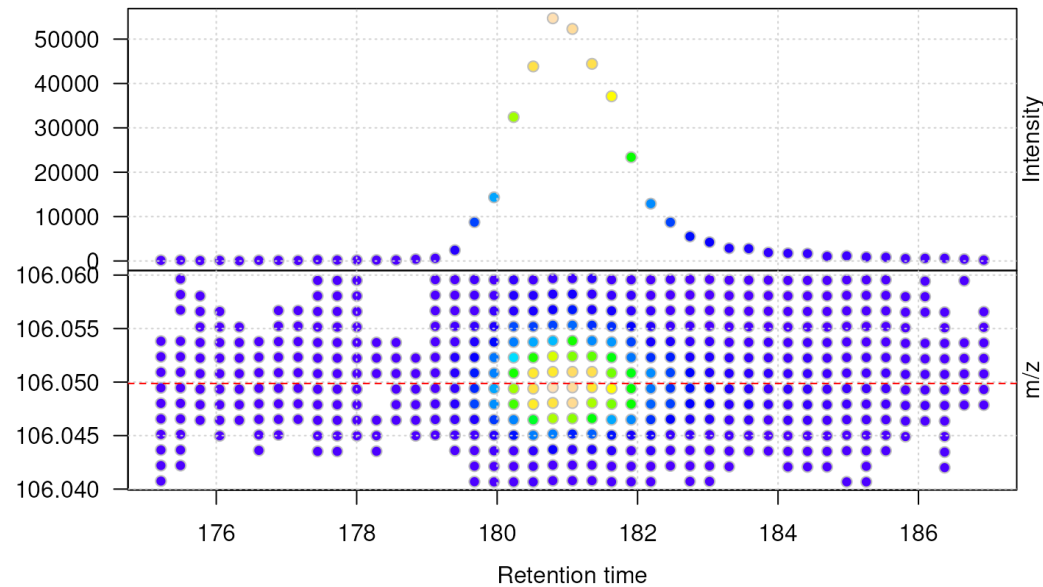
An Automated Pipeline for High-Throughput Label-Free Quantitative Proteomics

# Pre-processing

## Peak detection

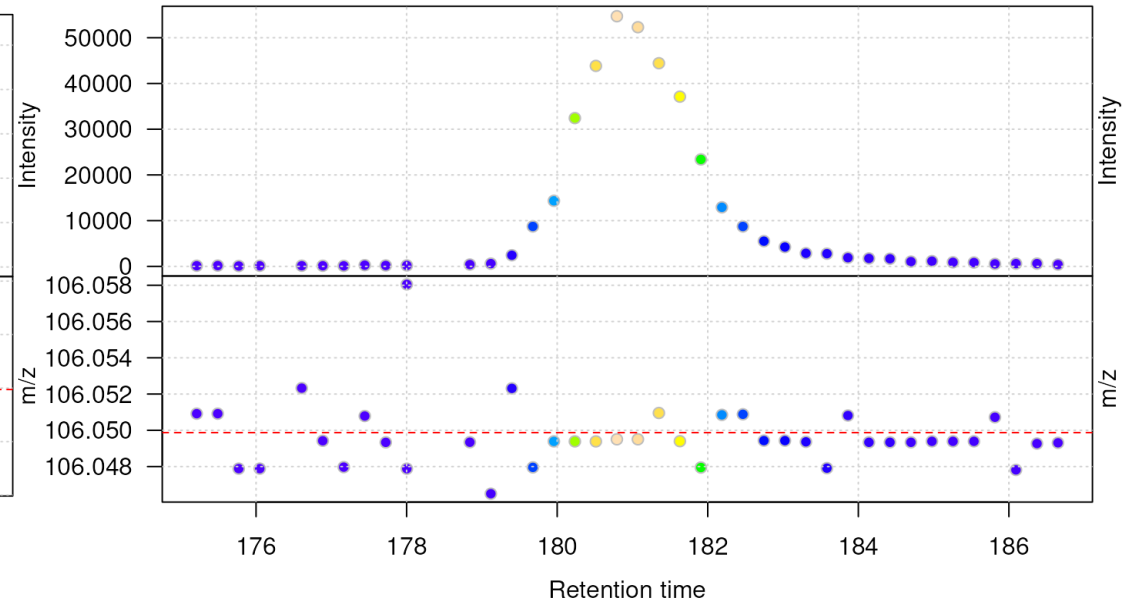
### Profile data

20171016\_POOL\_POS\_3\_105-134.mzML



### Centroid data

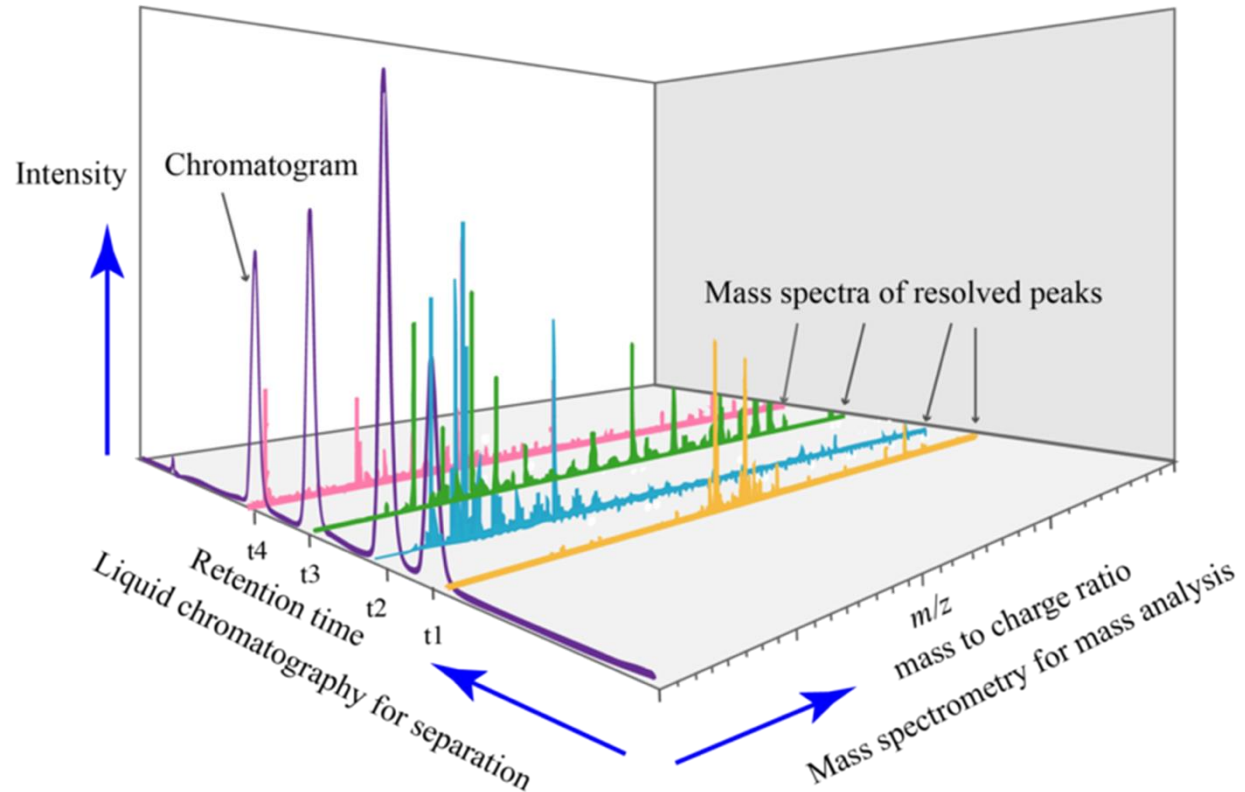
20171016\_POOL\_POS\_3\_105-134.mzML



The **goal** is to identify the **chromatographic peaks** (peaks across the retention time axis) and return a **peak feature** defined by m/z, rt and intensity per sample.

# Pre-processing

Peak detection with OpenMS FeatureFinderMetabo



1. Mass trace detection  
`MassTraceDetection()`
2. Elution peak detection  
`ElutionPeakDetection()`
3. Feature assembly  
`FeatureFindingMetabo()`

Kenar, E et al. (2014). Automated label-free quantification of metabolites from liquid chromatography-mass spectrometry data. *Molecular & cellular proteomics*.

# Pre-processing

Peak detection with OpenMS FeatureFinderMetabo

## 1. Mass trace detection

$P = \{p_k\}$  where  $p_k = \{rt_k, m/z_k, I_k\}$

Scan	RT	m/z	Intensity
1	240	181.0710	100
2	241	181.0711	200
3	242	181.0709	500
4	243	181.0710	300
1	240	455.2802	80
2	241	455.2800	120
3	242	455.2803	300

# Pre-processing

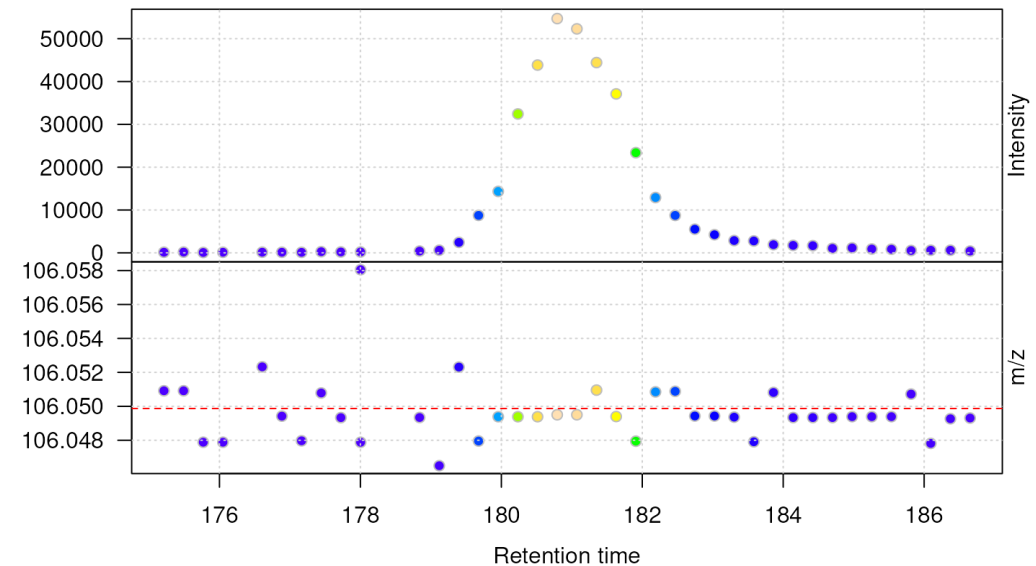
Peak detection with OpenMS FeatureFinderMetabo

## 1. Mass trace detection

$P = \{p_k\}$  where  $p_k = \{rt_k, m/z_k, I_k\}$

Scan	RT	m/z	Intensity
1	240	181.0710	100
2	241	181.0711	200
3	242	181.0709	500
4	243	181.0710	300
1	240	455.2802	80
2	241	455.2800	120
3	242	455.2803	300

Mass trace T: List of  $n$  peaks  $p_k \in P$  with similar  $m/z$  and that occur in adjacent scans of an LC-MS run (they elute at similar RT).



# Pre-processing

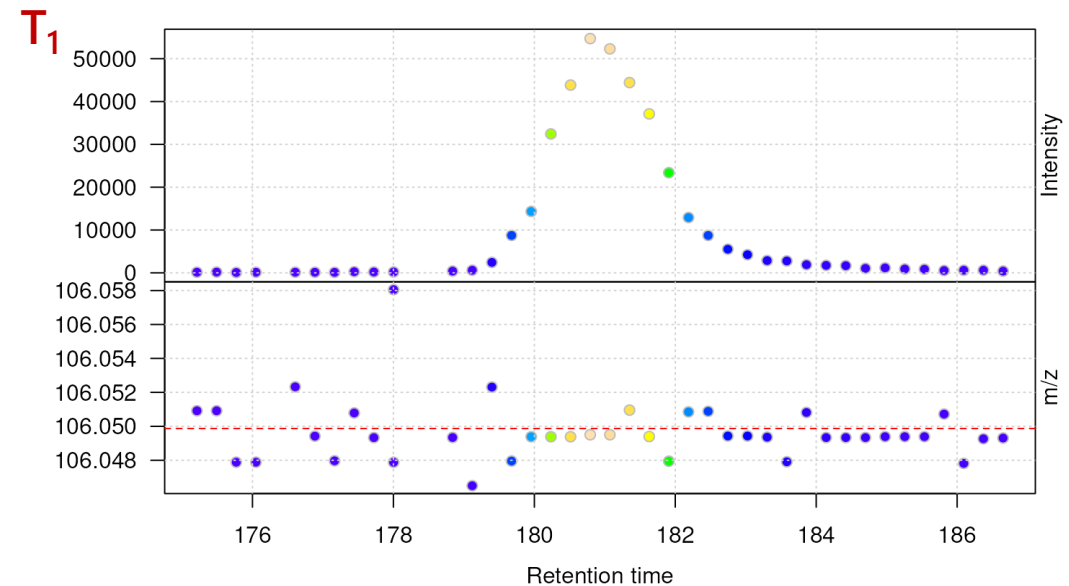
Peak detection with OpenMS FeatureFinderMetabo

## 1. Mass trace detection

$P = \{p_k\}$  where  $p_k = \{rt_k, m/z_k, I_k\}$

Scan	RT	m/z	Intensity
1	240	181.0710	100
2	241	181.0711	200
3	242	181.0709	500
4	243	181.0710	300
1	240	455.2802	80
2	241	455.2800	120
3	242	455.2803	300

Mass trace T: List of  $n$  peaks  $p_k \in P$  with similar  $m/z$  and that occur in adjacent scans of an LC-MS run (they elute at similar RT).



# Pre-processing

Peak detection with OpenMS FeatureFinderMetabo

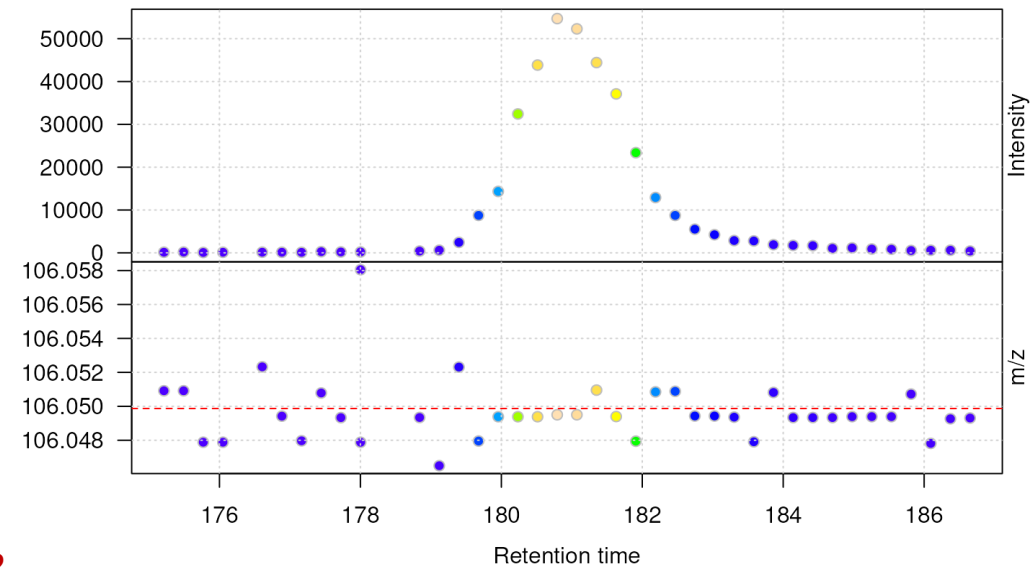
## 1. Mass trace detection

$P = \{p_k\}$  where  $p_k = \{rt_k, m/z_k, I_k\}$

Scan	RT	m/z	Intensity
1	240	181.0710	100
2	241	181.0711	200
3	242	181.0709	500
4	243	181.0710	300
1	240	455.2802	80
2	241	455.2800	120
3	242	455.2803	300

$T_2$

Mass trace  $T$ : List of  $n$  peaks  $p_k \in P$  with similar  $m/z$  and that occur in adjacent scans of an LC-MS run (they elute at similar RT).



# Pre-processing – Peak detection

OpenMS FeatureFinderMetabo

## 2. Elution peak detection

Does this mass trace contain multiple chromatographic peaks?

# Pre-processing – Peak detection

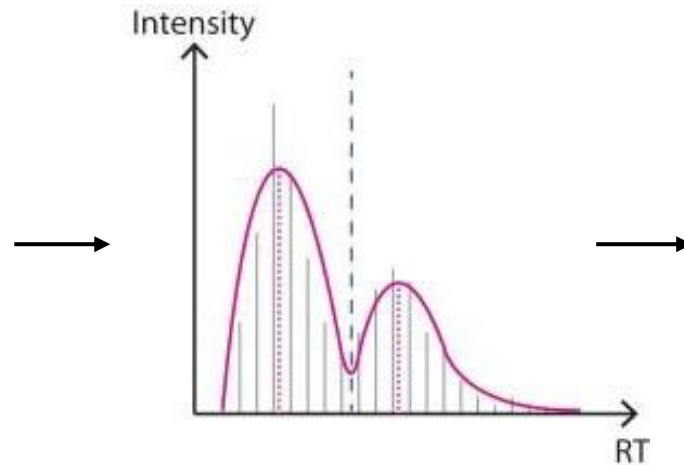
OpenMS FeatureFinderMetabo

## 2. Elution peak detection

Does this mass trace contain multiple chromatographic peaks?

Two compounds can have:

- Almost identical  $m/z$
- Overlap in RT



1. Smooth the signal with LOWESS
2. Find:
  - Local maxima (peaks)
  - Minima between maxima (separation points)
3. Each split becomes a final mass trace

**Intensity:** Area under the chromatographic peak

**Retention time:** Position of the chromatographic peak maximum

# Pre-processing – Peak detection

OpenMS FeatureFinderMetabo

## 3. Feature assembly

Which traces belong together as one metabolite feature  
(find isotopic patterns)?

# Pre-processing – Peak detection

OpenMS FeatureFinderMetabo

## 3. Feature assembly

Which traces belong together as one metabolite feature (find isotopic patterns)?

Scan	RT	m/z	Intensity
1	240	181.0710	100
2	241	181.0711	200
3	242	181.0709	500
4	243	181.0710	300
1	240	455.2802	80
2	241	455.2800	120
3	242	455.2803	300



Trace	RT	m/z
T <sub>1</sub>	241.85	181.0710
T <sub>2</sub>	243	455.2800

# Pre-processing – Peak detection

OpenMS FeatureFinderMetabo

## 3. Feature assembly

Which traces belong together as one metabolite feature (find isotopic patterns)?

Trace	RT	m/z
T <sub>1</sub>	241.85	181.0710
T <sub>3</sub>	242.75	182.0740
T <sub>4</sub>	242.02	183.0770
T <sub>5</sub>	242.52	184.0800

# Pre-processing – Peak detection

OpenMS FeatureFinderMetabo

## 3. Feature assembly

Which traces belong together as one metabolite feature (find isotopic patterns)?

Trace	RT	m/z
T <sub>1</sub>	241.85	181.0710
T <sub>3</sub>	242.75	182.0740
T <sub>4</sub>	242.02	183.0770
T <sub>5</sub>	242.52	184.0800

<sup>13</sup>C vs <sup>12</sup>C

The diagram shows four curved arrows pointing from the m/z values of T<sub>1</sub>, T<sub>3</sub>, T<sub>4</sub>, and T<sub>5</sub> to the right. Each arrow is labeled with '+1.003', indicating the mass difference between adjacent isotopic peaks.

# Pre-processing – Peak detection

OpenMS FeatureFinderMetabo

## 3. Feature assembly

Which traces belong together as one metabolite feature (find isotopic patterns)?

Trace	RT	m/z
T <sub>1</sub>	241.85	181.0710
T <sub>3</sub>	242.75	182.0740
T <sub>4</sub>	242.02	183.0770
T <sub>5</sub>	242.52	184.0800

<sup>13</sup>C vs <sup>12</sup>C

+1.003  
+1.003  
+1.003



Mass traces co-elute  
Correct m/z distances  
Correct isotopic abundance ratios

# Pre-processing – Peak detection

OpenMS FeatureFinderMetabo

## 3. Feature assembly

Which traces belong together as one metabolite feature (find isotopic patterns)?

Trace	RT	m/z
T <sub>1</sub>	241.85	181.0710
T <sub>3</sub>	242.75	182.0740
T <sub>4</sub>	242.02	183.0770
T <sub>5</sub>	242.52	184.0800

<sup>13</sup>C vs <sup>12</sup>C

+1.003  
+1.003  
+1.003

→ Feature 1



Mass traces co-elute  
Correct m/z distances  
Correct isotopic abundance ratios

# Pre-processing – Peak detection

OpenMS FeatureFinderMetabo

## 3. Feature assembly

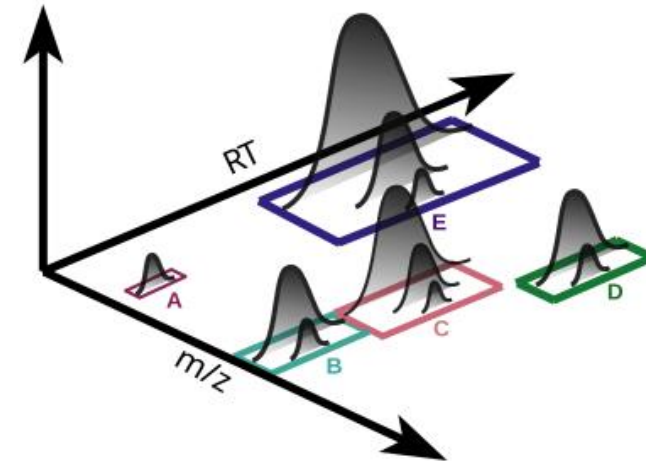
Which traces belong together as one metabolite feature (find isotopic patterns)?

Trace	RT	m/z
T <sub>1</sub>	241.85	181.0710
T <sub>3</sub>	242.75	182.0740
T <sub>4</sub>	242.02	183.0770
T <sub>5</sub>	242.52	184.0800

<sup>13</sup>C vs <sup>12</sup>C

+1.003  
+1.003  
+1.003

→ Feature 1



Mass traces co-elute  
Correct m/z distances  
Correct isotopic abundance ratios

# Pre-processing – Peak detection

OpenMS FeatureFinderMetabo

What do we have now?

Sample 1

m/z	RT	Intensity
m/z <sub>1</sub>	RT <sub>1</sub>	I <sub>1</sub>
m/z <sub>2</sub>	RT <sub>2</sub>	I <sub>2</sub>
...	...	...
m/z <sub>n</sub>	RT <sub>n</sub>	I <sub>n</sub>

Sample 2

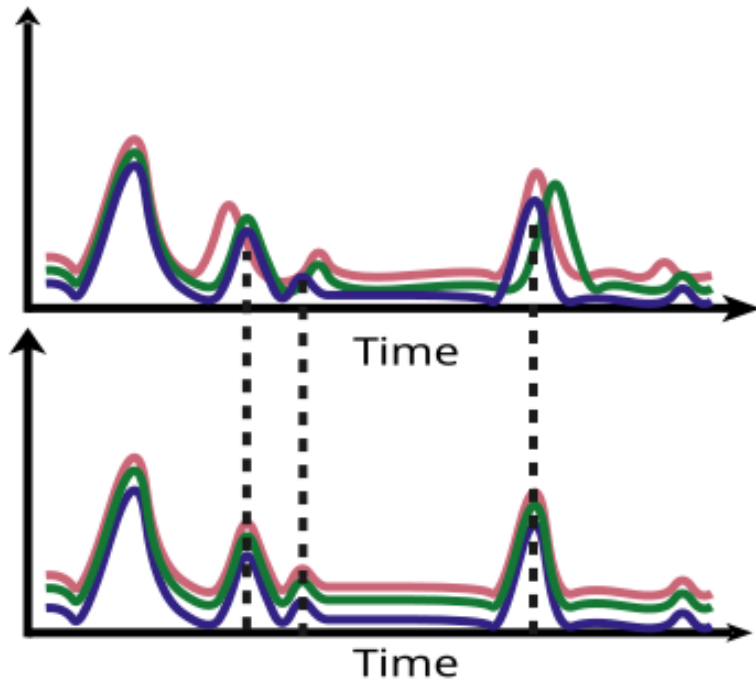
m/z	RT	Intensity
m/z <sub>1</sub>	RT <sub>1</sub>	I <sub>1</sub>
m/z <sub>2</sub>	RT <sub>2</sub>	I <sub>2</sub>
...	...	...
m/z <sub>t</sub>	RT <sub>t</sub>	I <sub>t</sub>

Sample k

m/z	RT	Intensity
m/z <sub>1</sub>	RT <sub>1</sub>	I <sub>1</sub>
m/z <sub>2</sub>	RT <sub>2</sub>	I <sub>2</sub>
...	...	...
m/z <sub>m</sub>	RT <sub>m</sub>	I <sub>m</sub>

# Pre-processing – Peak alignment and matching

Goal: Match features across different samples together



1. Peak alignment across the retention time axis

`MapAlignmentAlgorithmPoseClustering()`

2. Grouping corresponding features in multiple runs

`FeatureGroupingAlgorithmKD()`



**Break!**

# Pre-processing – Hands on

# Click on: Open in GitHub Codespaces

Mass spectrometry-based  
Metabolomics introduction  
documentation



Search

⌘ + K

Metaboigniter

MS-based metabolomics

Metaboigniter



## Metaboigniter

### Open GitHub codespace

Use the following link to open a GitHub codespace with most of the required software installed:

⚠ If you do it manually, make sure to select the bigger machine with 4 cores and 16GB RAM



Open in GitHub Codespaces

# Alternative setup of codespace

1

2

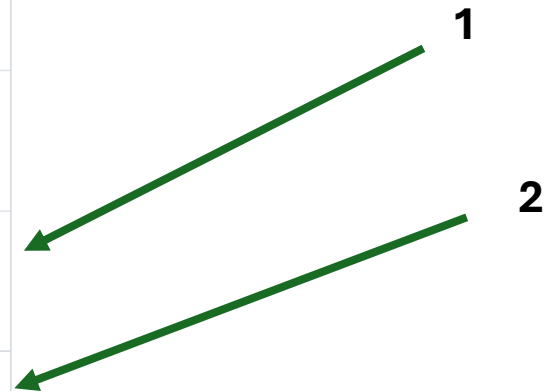
The screenshot shows the GitHub interface for the repository `biosustain / dsp_course_metabolomics_intro`. A blue notification banner at the top states "Your main branch isn't protected". Below this, the repository navigation bar includes "Code", "Issues", "Pull requests", "Agents", "Actions", "Projects", "Wiki", "Security and quality", "Insights", and "Settings". The repository name and "Public" status are displayed, along with "Edit Pins", "Watch 0", "Fork 1", and "Star" buttons. A "Code" button is highlighted with a green arrow. A dropdown menu is open from the "Code" button, showing options for "Local" and "Codespaces". The "Codespaces" section includes a plus sign, a three-dot menu (highlighted with a green arrow), and a list of actions: "New with options...", "Configure dev container", "Set up prebuilds", "Manage codespaces", "Share a deep link", and "What are codespaces?". The repository file list is visible on the left, and the "About" section is on the right.

# Codespace with 4 cores (and 16GB of memory)

- **Not Safari**
- **Ignore Errors (todo tree)**

Create codespace for  
**biosustain/dsp\_course\_metabolomics\_intro**

<b>Branch</b> This branch will be checked out on creation	<input type="text" value="main"/>
<b>Dev container configuration</b> Your codespace will use this configuration	<input type="text" value="nextflow-training"/>
<b>Region</b> Your codespace will run in the selected region	<input type="text" value="Europe West"/>
<b>Machine type</b> Resources for your codespace	<input type="text" value="4-core"/>
<input type="button" value="Create codespace"/>	



**Your main branch isn't protected** Dismiss Protect this branch

Protect this branch from force pushing or deletion, or require status checks before merging. [View documentation.](#)

main 4 Branches 0 Tags Go to file Add file Code

enryH and feliciaschulz Data analysis (#8) 6cd94bb · 48 minutes ago 9 Commits

.devcontainer	set latest dockerfile (#4)	4 days ago
.github/workflows	add docker image similar to proteomics course (#2)	last week
.vscode	set latest dockerfile (#4)	4 days ago
bin	Data analysis (#8)	48 minutes ago
data	Data analysis (#8)	48 minutes ago
material	Fix docs (#7)	11 hours ago

**About**

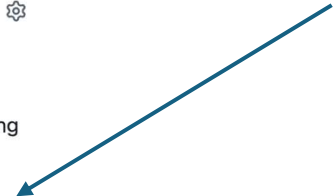
About Introduction to massspectrometry-based metabolomics (raw data processing and downstream analysis).

[biosustain.github.io/dsp\\_course\\_...](#)

- Readme
- Activity
- Custom properties
- 1 star
- 0 watching
- 1 fork

Report repository

Go to website



# Repo Overview

Feel free to ask questions at  
any point

enryH and feliciaschulz Data analysis (#8) <span>6cd94bb · 49 minutes ago</span> <span>🔄 9 Commits</span>	
📁 .devcontainer	🔗 set latest dockerfile (#4) 4 days ago
📁 .github/workflows	🔗 add docker image similar to proteomics course (#2) last week
📁 .vscode	🔗 set latest dockerfile (#4) 4 days ago
📁 bin	Data analysis (#8) 49 minutes ago
📁 data	Data analysis (#8) 49 minutes ago
📁 material	Fix docs (#7) 11 hours ago
📁 results_prepared	Data analysis (#8) 49 minutes ago
📄 .gitignore	🔗 add docker image similar to proteomics course (#2) last week
📄 2_data_analysis.ipynb	Data analysis (#8) 49 minutes ago
📄 2_data_analysis.py	Data analysis (#8) 49 minutes ago
📄 README.md	Updated Agenda and Location June 2026 (#6) 4 days ago
📄 conf.py	Data analysis (#8) 49 minutes ago
📄 index.md	Data analysis (#8) 49 minutes ago
📄 nextflow.config	🔗 set latest dockerfile (#4) 4 days ago
📄 requirements.txt	Data analysis (#8) 49 minutes ago
📄 requirements_docs.txt	🚀 initialize relatively bare bone repo for metabolomics cou... 3 months ago

# Nextflow and nf-core pipeline template



Nextflow is a workflow executor (analysis steps combined in an execution graph)

- Reproducible and scalable

Nf-core is a collection of workflows maintained by nextflow and an open-source community



A global community effort to collect a curated set of open-source analysis pipelines built using Nextflow.

## Pipelines

Browse the 130 pipelines that are currently available as part of nf-core.

Search... Released 80 Under development 38 Archived 12 Last release ▾ ☰ ☰

proteinfamilies ✓ ☆ 14  
Generation and update of protein families  
metagenomics protein-families proteomics  
1.1.0 released 4 days ago

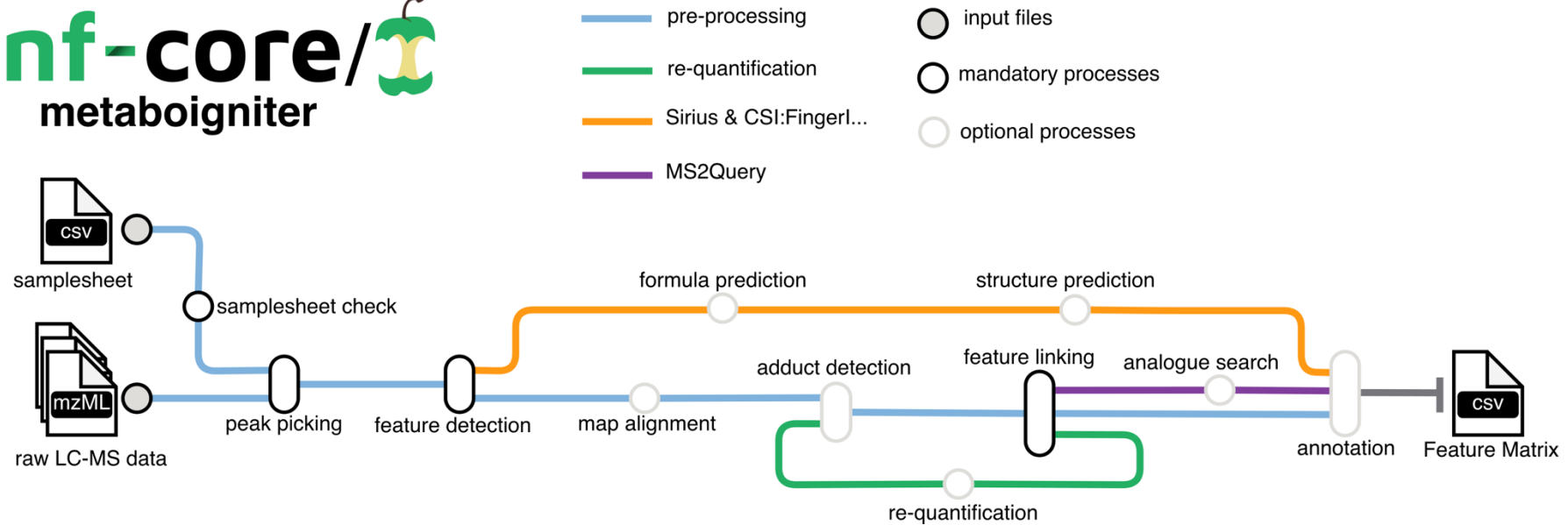
epitopeprediction ✓ ☆ 44  
A bioinformatics best-practice analysis pipeline for epitope prediction and annotation  
epitope epitope-prediction mhc-binding-prediction  
3.0.0 released 5 days ago

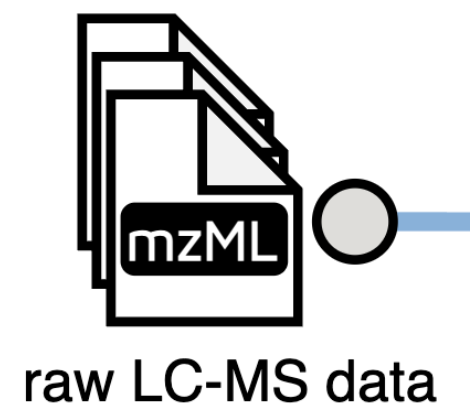
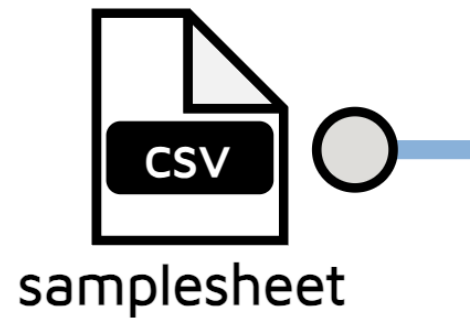
Works with

- Linux machines (on servers)
- on MacOS

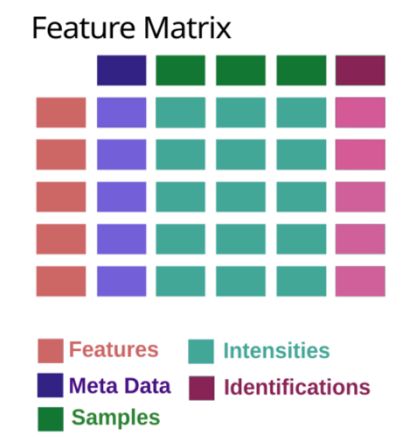
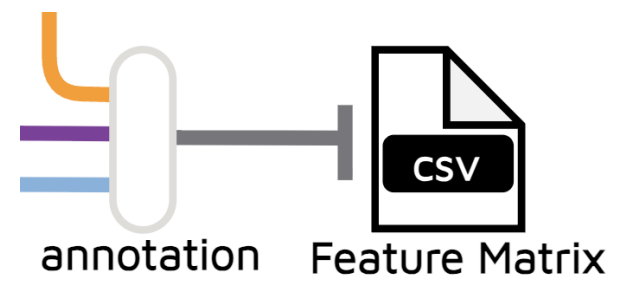
# INPUTS

# Metaboigniter – Metro Map





( ... )



# Samplesheet

sample	level	type	msfile
MS_A_POS	MS1	normal	data/MTBLS8735/MS_A_POS.mzML
MS_B_POS	MS1	normal	data/MTBLS8735/MS_B_POS.mzML
MS_C_POS	MS1	normal	data/MTBLS8735/MS_C_POS.mzML
MS_D_POS	MS1	normal	data/MTBLS8735/MS_D_POS.mzML
MS_E_POS	MS1	normal	data/MTBLS8735/MS_E_POS.mzML
MS_F_POS	MS1	normal	data/MTBLS8735/MS_F_POS.mzML
MS_QC_POOL_1_POS	MS1	normal	data/MTBLS8735/MS_QC_POOL_1_POS.mzML
MS_QC_POOL_2_POS	MS1	normal	data/MTBLS8735/MS_QC_POOL_2_POS.mzML
MS_QC_POOL_3_POS	MS1	normal	data/MTBLS8735/MS_QC_POOL_3_POS.mzML
MS_QC_POOL_4_POS	MS1	normal	data/MTBLS8735/MS_QC_POOL_4_POS.mzML
MSMS_2_A_CE20_POS	MS2	normal	data/MTBLS8735/MSMS_2_A_CE20_POS.mzML
MSMS_2_A_CE30_POS	MS2	normal	data/MTBLS8735/MSMS_2_A_CE30_POS.mzML
MSMS_2_A_CES_POS	MS2	normal	data/MTBLS8735/MSMS_2_A_CES_POS.mzML
MSMS_2_E_CE20_POS	MS2	normal	data/MTBLS8735/MSMS_2_E_CE20_POS.mzML
MSMS_2_E_CE30_POS	MS2	normal	data/MTBLS8735/MSMS_2_E_CE30_POS.mzML
MSMS_2_E_CES_POS	MS2	normal	data/MTBLS8735/MSMS_2_E_CES_POS.mzML

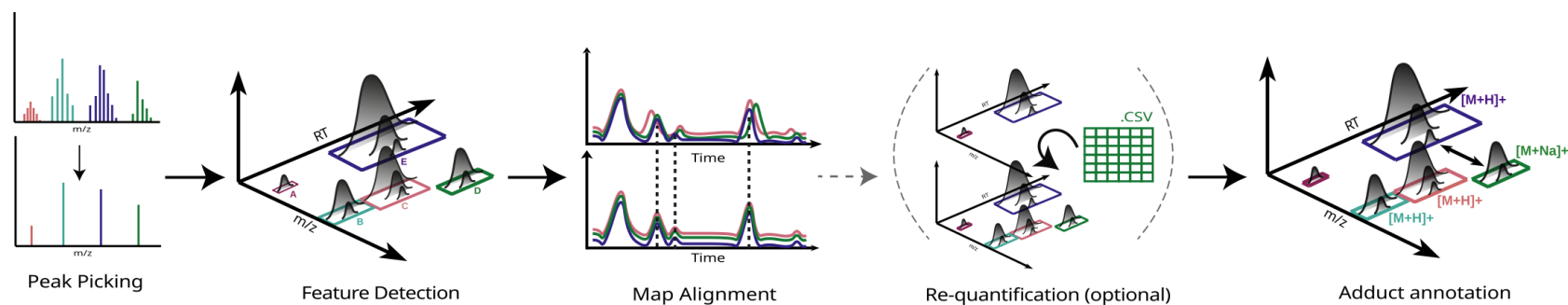
- Open standard, text based:
  - mzML files (indexed)
- Needs to be created from .d (bruker), .raw (Thermo), etc. vendor specific file formats
  - Conversion not (yet) integrated into workflow

# One MS1 spectrum (DDA)

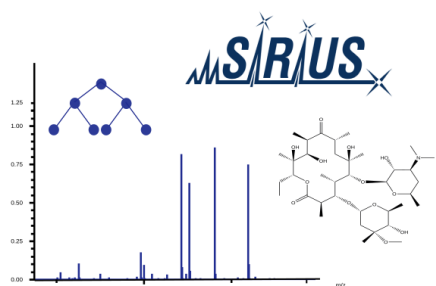
```
<spectrum id="controllerType=0 controllerNumber=1 scan=6" index="5" defaultArrayLength="736">
  <cvParam cvRef="MS" accession="MS:1000511" value="1" name="ms level" />
  <cvParam cvRef="MS" accession="MS:1000579" value="" name="MS1 spectrum" />
  <cvParam cvRef="MS" accession="MS:1000130" value="" name="positive scan" />
  <cvParam cvRef="MS" accession="MS:1000285" value="15240991" name="total ion current" />
  <cvParam cvRef="MS" accession="MS:1000127" value="" name="centroid spectrum" />
  <cvParam cvRef="MS" accession="MS:1000504" value="401.923461914063" name="base peak m/z" unitAccession="MS:1000040" unitName="m/z" unitCvRef="MS" />
  <cvParam cvRef="MS" accession="MS:1000505" value="2704116.25" name="base peak intensity" unitAccession="MS:1000131" unitName="number of detector counts" unitCvRef="MS" />
  <cvParam cvRef="MS" accession="MS:1000528" value="375.884124755859" name="lowest observed m/z" unitAccession="MS:1000040" unitName="m/z" unitCvRef="MS" />
  <cvParam cvRef="MS" accession="MS:1000527" value="1665.40612792969" name="highest observed m/z" unitAccession="MS:1000040" unitName="m/z" unitCvRef="MS" />
  <scanList count="1">
    <cvParam cvRef="MS" accession="MS:1000795" value="" name="no combination" />
    <scan instrumentConfigurationRef="IC1">
      <cvParam cvRef="MS" accession="MS:1000016" value="6.0323342" name="scan start time" unitAccession="U0:0000031" unitName="minute" unitCvRef="U0" />
      <cvParam cvRef="MS" accession="MS:1000512" value="FTMS + p NSI Full ms [375.0000-1800.0000]" name="filter string" />
      <cvParam cvRef="MS" accession="MS:1000927" value="50" name="ion injection time" unitAccession="U0:0000028" unitName="millisecond" unitCvRef="U0" />
      <scanWindowList count="1">
        <scanWindow>
          <cvParam cvRef="MS" accession="MS:1000501" value="375" name="scan window lower limit" unitAccession="MS:1000040" unitName="m/z" unitCvRef="MS" />
          <cvParam cvRef="MS" accession="MS:1000500" value="1800" name="scan window upper limit" unitAccession="MS:1000040" unitName="m/z" unitCvRef="MS" />
        </scanWindow>
      </scanWindowList>
    </scan>
  </scanList>
  <binaryDataArrayList count="2">
    <binaryDataArray encodedLength="3152">
      <cvParam cvRef="MS" accession="MS:1000514" value="" name="m/z array" unitAccession="MS:1000040" unitName="m/z" unitCvRef="MS" />
      <cvParam cvRef="MS" accession="MS:1000523" value="" name="64-bit float" />
      <cvParam cvRef="MS" accession="MS:1000574" value="" name="zlib compression" />
      <binary>eJwT03LYXUawPHjVumjGA1LG5dTGtMi3MYUG32496hJOiNX3EvLcJLNsEDNckGWIyoGxaIwiU3AsUwcZ2TJEm5Wimg48C466j15M12ZRgU1paCTPP+/399Xne9fceFk3TzIeylxna//0Vjck5ovM26h1oDLdI
    </binaryDataArray>
    <binaryDataArray encodedLength="4356">
      <cvParam cvRef="MS" accession="MS:1000515" value="" name="intensity array" unitAccession="MS:1000131" unitName="number of counts" unitCvRef="MS" />
      <cvParam cvRef="MS" accession="MS:1000523" value="" name="64-bit float" />
      <cvParam cvRef="MS" accession="MS:1000574" value="" name="zlib compression" />
      <binary>eJwTWHlCt+kXviLFSI2MTJi5ZWIaTTGdSl+EkT3KVtypUVl+dhR5TG4qo7Jk8FPWm32pjC1kuxUleU0JuSVDPozK0pdKpj3P/PV+3vu+73nPec45zZnvFQRBUrdyQxIEQWzvf6ppFBz63sIYceZq06jafJHbNGr;
    </binaryDataArray>
  </binaryDataArrayList>
</spectrum>
```

# Metaboigniter - Overview

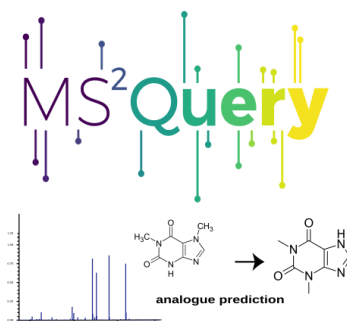
## Pre-Processing



## Identification (optional)

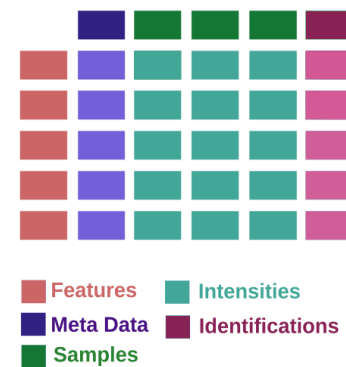


Predict metabolite formulas and structures with SIRIUS and CSI:FingerID.



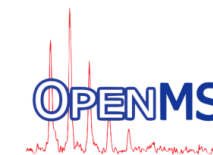
Detect analogue molecules from your MS2 spectra with the machine learning tool MS2Query.

## Feature Matrix




Feature Linking

pyOpenMS



# OpenMS and pyOpenMS



OpenMS 3.5.0 documentation

Search

ABOUT

- Installation
- Community
- Learning

GETTING STARTED

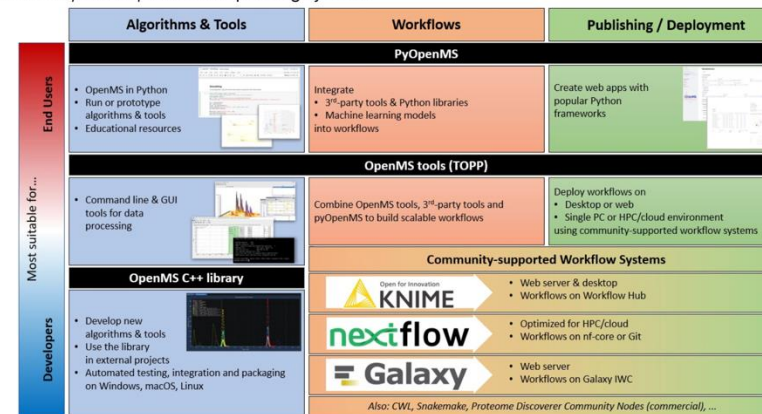
- Introduction
- WebApps
- Workflow Editor
- TOPPView
- TOPP Tools
- pyOpenMS

MANUAL

- Contribute
- Developers
- Additional

## 1. What is OpenMS

OpenMS is a free, open-source framework based on a C++ library with Python bindings. It is commonly used for liquid chromatography-mass spectrometry (LC-MS) data management and analyses. OpenMS provides an infrastructure for the rapid development of mass spectrometry related software as well as a rich toolset built on top of it. OpenMS is available under the [three clause BSD licence](#) and runs under Windows, macOS, and Linux operating systems.



OpenMS developers can create new C++ algorithms and tools, while users can execute tools or implement new algorithms or scripts in Python. Workflows integrate pyOpenMS scripts and OpenMS tools with third-party tools and external Python libraries to create scalable data-processing pipelines. For deployment, users can use pyOpenMS with web frameworks or deploy workflows on desktop, high-performance computing (HPC) or cloud infrastructure using one of the community-supported workflow systems.

OpenMS supports the Proteomics Standard Initiative (PSI) formats for MS data. The main contributors of OpenMS are currently the Eberhard-Karls-Universität in Tübingen, the Freie Universität Berlin, and the University of Toronto.

# pyOpenMS functions highlighted

Mass spectrometry-based  
Metabolomics introduction  
documentation

🔍

🔍 Search ✖ + K

**Metaboigniter**

- MS-based metabolomics
- Metaboigniter
- Preprocessing manually using pyOpenMS**

**Data Analysis**

- Downstream Data Analysis
- Hands-On

**Interpretation**

- Hands-On

**Datasets**

- Datasets overview
- Data

**Codespace**

- Container info

☰

🔍 🔄 ⬇️ 🗄️ 🌐

## Data visualization

Before starting the pre-processing workflow, we visualize the **Base Peak Chromatogram (BPC)** for each sample. The BPC represents the intensity of the most intense ion detected in every MS1 spectrum over the course of the chromatographic run. Plotting the BPC provides a quick overview of the data quality and allows us to compare the overall signal profiles between samples. Similar chromatographic patterns across runs suggest consistent instrument performance, while large differences may indicate potential issues that should be investigated before proceeding with further analysis.

▶ Show code cell source

The figure is a line plot titled "BPC" showing the Base Peak Chromatogram. The y-axis is labeled "Intensity" and has a multiplier of  $1e6$ , with values ranging from 0.0 to 1.2. The x-axis is labeled "Retention Time (s)" and ranges from 0 to 500. The plot displays three data series: QC (orange), CVD (blue), and CTR (purple). Each series shows a series of peaks, with the most prominent peaks occurring between 0 and 200 seconds. The QC series has the highest intensity peaks, reaching approximately 1.2 at around 20 and 130 seconds. The CVD and CTR series show lower intensity peaks, with CVD reaching about 1.1 at 130 seconds and CTR reaching about 0.4 at 130 seconds. The plot also shows a legend on the right side with the following entries: QC (orange), CVD (blue), CTR (purple), QC (orange), CTR (purple), CVD (blue), QC (orange), CTR (purple), CVD (blue), and QC (orange).

# Data analysis - Filtering and normalization

Input

Peak	m/z	rt	$I_1$	...	$I_k$
$P_1$	$m/z_1$	$rt_1$	$I_{11}$	...	$I_{1k}$
...	...	...	...	...	...
$P_n$	$m/z_n$	$rt_n$	$I_{n1}$	...	$I_{nk}$

# Data analysis - Filtering and normalization

Input

Peak	m/z	rt	$I_1$	...	$I_k$
$P_1$	$m/z_1$	$rt_1$	$I_{11}$	...	$I_{1k}$
...	...	...	...	...	...
$P_n$	$m/z_n$	$rt_n$	$I_{n1}$	...	$I_{nk}$

# Data analysis - Filtering and normalization

Input

Peak	m/z	rt	$I_1$	...	$I_k$
$P_1$	$m/z_1$	$rt_1$	$I_{11}$	...	$I_{1k}$
...	...	...	...	...	...
$P_n$	$m/z_n$	$rt_n$	$I_{n1}$	...	$I_{nk}$

# Data analysis - Filtering and normalization

Input

Peak	m/z	rt	$I_1$	...	$I_k$
$P_1$	$m/z_1$	$rt_1$	$I_{11}$	...	$I_{1k}$
...	...	...	...	...	...
$P_n$	$m/z_n$	$rt_n$	$I_{n1}$	...	$I_{nk}$

# Data analysis - Filtering and normalization


Input

Peak	m/z	rt	$I_1$	...	$I_k$
$P_1$	$m/z_1$	$rt_1$	$I_{11}$	...	$I_{1k}$
...	...	...	...	...	...
$P_n$	$m/z_n$	$rt_n$	$I_{n1}$	...	$I_{nk}$

# Data analysis - Filtering and normalization

Input

Peak	m/z	rt	$I_1$	...	$I_k$
$P_1$	$m/z_1$	$rt_1$	$I_{11}$	...	$I_{1k}$
...	...	...	...	...	...
$P_n$	$m/z_n$	$rt_n$	$I_{n1}$	...	$I_{nk}$



Missingness due to errors or real absence of the molecule

# Data analysis - Filtering and normalization

Filtering using the 80% rule

Peak	m/z	rt	missingness (%)	$l_1$	...	$l_k$
$P_1$	$m/z_1$	$rt_1$	5%	$l_{11}$	...	$l_{1k}$
$P_2$	$m/z_2$	$rt_2$	8%	$l_{21}$	...	$l_{2k}$
$P_3$	$m/z_3$	$rt_3$	22%	$l_{31}$	...	$l_{3k}$
...	...	...	...	...	...	...
$P_n$	$m/z_n$	$rt_n$	42%	$l_{n1}$	...	$l_{nk}$

# Data analysis - Filtering and normalization

Filtering using the 80% rules

Peak	m/z	rt	missingness (%)	$l_1$	...	$l_k$
$P_1$	$m/z_1$	$rt_1$	5%	$l_{11}$	...	$l_{1k}$
$P_2$	$m/z_2$	$rt_2$	8%	$l_{21}$	...	$l_{2k}$
<del><math>P_3</math></del>	<del><math>m/z_3</math></del>	<del><math>rt_3</math></del>	<del>22%</del>	<del><math>l_{31}</math></del>	<del>...</del>	<del><math>l_{3k}</math></del>
...	...	...	...	...	...	...
<del><math>P_n</math></del>	<del><math>m/z_n</math></del>	<del><math>rt_n</math></del>	<del>42%</del>	<del><math>l_{n1}</math></del>	<del>...</del>	<del><math>l_{nk}</math></del>

# Data analysis - Filtering and normalization

Filtering using the 80% rules

Peak	m/z	rt	missingness (%)	$l_1$	...	$l_k$
$P_1$	$m/z_1$	$rt_1$	5%	$l_{11}$	...	$l_{1k}$
$P_2$	$m/z_2$	$rt_2$	8%	$l_{21}$	...	$l_{2k}$
<del><math>P_3</math></del>	<del><math>m/z_3</math></del>	<del><math>rt_3</math></del>	<del>22%</del>	<del><math>l_{31}</math></del>	<del>...</del>	<del><math>l_{3k}</math></del>
...	...	...	...	...	...	...
<del><math>P_n</math></del>	<del><math>m/z_n</math></del>	<del><math>rt_n</math></del>	<del>42%</del>	<del><math>l_{n1}</math></del>	<del>...</del>	<del><math>l_{nk}</math></del>



Positive acquisition mode: 2297 peaks

Negative acquisition mode: 2063 peaks

# Data analysis - Filtering and normalization

Filtering using the modified 80% rule

Peak	m/z	rt	D <sub>1</sub>	D <sub>2</sub>	D <sub>k</sub>	D:missign ess (%)	C <sub>1</sub>	C <sub>2</sub>	C <sub>k</sub>	C:missig ness (%)	...
P <sub>1</sub>	m/z <sub>1</sub>	rt <sub>1</sub>	D <sub>11</sub>	D <sub>12</sub>	D <sub>1k</sub>	5%	C <sub>11</sub>	C <sub>12</sub>	C <sub>1k</sub>	7%	...
P <sub>2</sub>	m/z <sub>2</sub>	rt <sub>2</sub>	D <sub>21</sub>	D <sub>22</sub>	D <sub>2k</sub>	2%	C <sub>21</sub>	C <sub>22</sub>	C <sub>2k</sub>	83%	...
P <sub>3</sub>	m/z <sub>3</sub>	rt <sub>3</sub>	D <sub>31</sub>	D <sub>32</sub>	D <sub>3k</sub>	22%	C <sub>31</sub>	C <sub>32</sub>	C <sub>3k</sub>	18%	...
...	...	...	...	...	...	...	...	...	...	...	...
<del>P<sub>n</sub></del>	<del>m/z<sub>n</sub></del>	<del>rt<sub>n</sub></del>	<del>D<sub>41</sub></del>	<del>D<sub>42</sub></del>	<del>D<sub>4k</sub></del>	<del>42%</del>	<del>C<sub>41</sub></del>	<del>C<sub>42</sub></del>	<del>C<sub>4k</sub></del>	<del>21%</del>	<del>...</del>

# Data analysis - Filtering and normalization

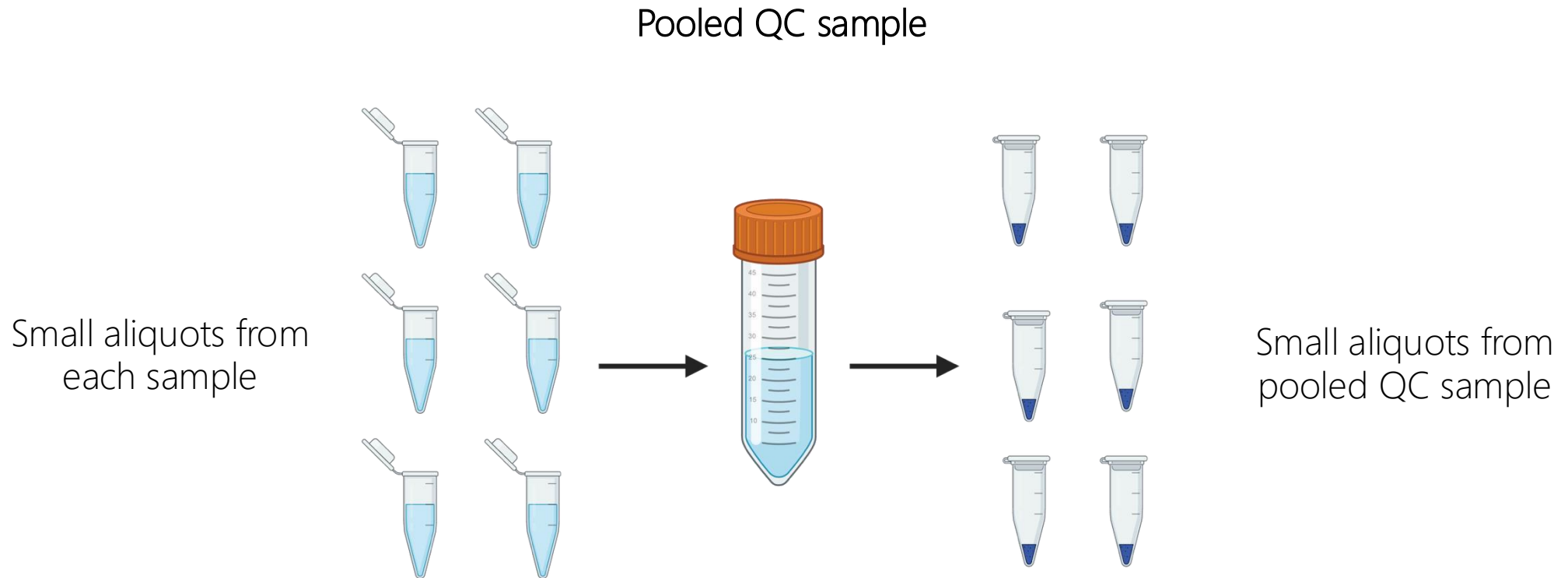
Filtering using the modified 80% rule

Peak	m/z	rt	D <sub>1</sub>	D <sub>2</sub>	D <sub>k</sub>	D:missign ess (%)	C <sub>1</sub>	C <sub>2</sub>	C <sub>k</sub>	C:missig ness (%)	...
P <sub>1</sub>	m/z <sub>1</sub>	rt <sub>1</sub>	D <sub>11</sub>	D <sub>12</sub>	D <sub>1k</sub>	5%	C <sub>11</sub>	C <sub>12</sub>	C <sub>1k</sub>	7%	...
P <sub>2</sub>	m/z <sub>2</sub>	rt <sub>2</sub>	D <sub>21</sub>	D <sub>22</sub>	D <sub>2k</sub>	2%	C <sub>21</sub>	C <sub>22</sub>	C <sub>2k</sub>	83%	...
P <sub>3</sub>	m/z <sub>3</sub>	rt <sub>3</sub>	D <sub>31</sub>	D <sub>32</sub>	D <sub>3k</sub>	22%	C <sub>31</sub>	C <sub>32</sub>	C <sub>3k</sub>	18%	...
...	...	...	...	...	...	...	...	...	...	...	...
<del>P<sub>n</sub></del>	<del>m/z<sub>n</sub></del>	<del>rt<sub>n</sub></del>	<del>D<sub>41</sub></del>	<del>D<sub>42</sub></del>	<del>D<sub>4k</sub></del>	<del>42%</del>	<del>C<sub>41</sub></del>	<del>C<sub>42</sub></del>	<del>C<sub>4k</sub></del>	<del>21%</del>	<del>...</del>

“Perfect biomarker” gets retained

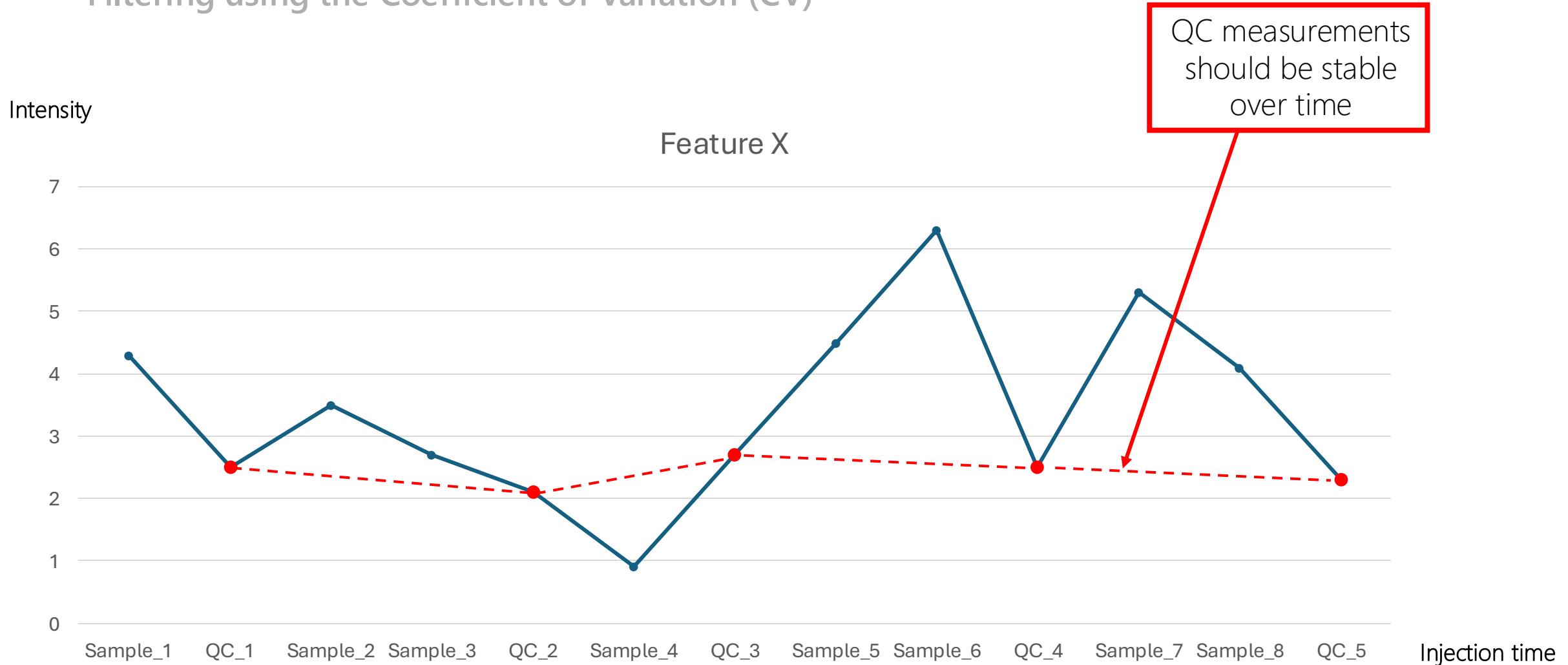
# Data analysis - Filtering and normalization

Filtering using the Coefficient of Variation (CV)



# Data analysis - Filtering and normalization

Filtering using the Coefficient of Variation (CV)



# Data analysis - Filtering and normalization

Filtering using the Coefficient of Variation (CV)

$$CV (\%) = \left( \frac{\text{Standard deviation}}{\text{Mean}} \right) \times 100$$

- If the QC CV > S CV: remove feature
- If the QC CV > threshold %: remove feature

Peak	m/z	rt	S <sub>1</sub>	S <sub>2</sub>	S <sub>k</sub>	S:CV	QC <sub>1</sub>	QC <sub>2</sub>	QC <sub>k</sub>	QC:CV	...
P <sub>1</sub>	m/z <sub>1</sub>	rt <sub>1</sub>	S <sub>11</sub>	S <sub>12</sub>	S <sub>1k</sub>	8%	QC <sub>11</sub>	QC <sub>12</sub>	QC <sub>1k</sub>	7%	...
P <sub>2</sub>	m/z <sub>2</sub>	rt <sub>2</sub>	S <sub>21</sub>	S <sub>22</sub>	S <sub>2k</sub>	10%	QC <sub>21</sub>	QC <sub>22</sub>	QC <sub>2k</sub>	5%	...
<del>P<sub>3</sub></del>	<del>m/z<sub>3</sub></del>	<del>rt<sub>3</sub></del>	<del>S<sub>31</sub></del>	<del>S<sub>32</sub></del>	<del>S<sub>3k</sub></del>	<del>5%</del>	<del>QC<sub>31</sub></del>	<del>QC<sub>32</sub></del>	<del>QC<sub>3k</sub></del>	<del>20%</del>	<del>...</del>
...	...	...	...	...	...	...	...	...	...	...	...
<del>P<sub>n</sub></del>	<del>m/z<sub>n</sub></del>	<del>rt<sub>n</sub></del>	<del>S<sub>41</sub></del>	<del>S<sub>42</sub></del>	<del>S<sub>4k</sub></del>	<del>12%</del>	<del>QC<sub>41</sub></del>	<del>QC<sub>42</sub></del>	<del>QC<sub>4k</sub></del>	<del>35%</del>	<del>...</del>

# Data analysis - Filtering and normalization

## Filtering using the Blanks control

### Blank extracts

- Prepared with same materials, chemicals, consumables as study samples
- But: no biological sample



### Why? Check for...

- ... compounds from other sources than tested biological material
- ... carryover in instrument

# Data analysis - Filtering and normalization

## Filtering using the Blanks control

Blank extracts

- Prepared with same materials, chemicals, consumables as study samples
- But: no biological sample

If the ratio of **total blanks intensities / total samples intensities** for a feature exceeds a certain threshold, remove

Peak	m/z	rt	S <sub>1</sub>	S <sub>2</sub>	S <sub>k</sub>	S:Total	B	B <sub>2</sub>	B <sub>k</sub>	B:Total	B_Total/ S_Total
P <sub>1</sub>	m/z <sub>1</sub>	rt <sub>1</sub>	S <sub>11</sub>	S <sub>12</sub>	S <sub>1k</sub>	...	B <sub>11</sub>	B <sub>12</sub>	B <sub>1k</sub>	...	0.2
<del>P<sub>2</sub></del>	<del>m/z<sub>2</sub></del>	<del>rt<sub>2</sub></del>	<del>S<sub>21</sub></del>	<del>S<sub>22</sub></del>	<del>S<sub>2k</sub></del>	<del>...</del>	<del>B<sub>21</sub></del>	<del>B<sub>22</sub></del>	<del>B<sub>2k</sub></del>	<del>...</del>	<del>1</del>
P <sub>3</sub>	m/z <sub>3</sub>	rt <sub>3</sub>	S <sub>31</sub>	S <sub>32</sub>	S <sub>3k</sub>	...	B <sub>31</sub>	B <sub>32</sub>	B <sub>3k</sub>	...	0.001
...	...	...	...	...	...	...	...	...	...	...	...
<del>P<sub>n</sub></del>	<del>m/z<sub>n</sub></del>	<del>rt<sub>n</sub></del>	<del>S<sub>41</sub></del>	<del>S<sub>42</sub></del>	<del>S<sub>4k</sub></del>	<del>...</del>	<del>B<sub>41</sub></del>	<del>B<sub>42</sub></del>	<del>B<sub>4k</sub></del>	<del>...</del>	<del>2</del>

# Data analysis - Imputation

## Imputation with zeros

After filtering: still some missing values → need imputation

Peak	m/z	rt	missingness (%)	$l_1$	$l_2$	$l_3$	...	$l_k$
$P_1$	$m/z_1$	$rt_1$	5%	$l_{11}$	N/A	$l_{13}$	...	$l_{1k}$
$P_2$	$m/z_2$	$rt_2$	8%	$l_{21}$	$l_{22}$	N/A	...	$l_{2k}$
$P_3$	$m/z_3$	$rt_3$	3%	N/A	$l_{32}$	$l_{33}$	...	$l_{3k}$
...	...	...	...	...	...	...	...	...
$P_n$	$m/z_n$	$rt_n$	2%	$l_{n1}$	$l_{n2}$	$l_{n3}$	...	$l_{nk}$

# Data analysis - Imputation

## Imputation with zeros

After filtering: still some missing values → need imputation

Common imputation strategy: Imputation with **zeros**

Peak	m/z	rt	missingness (%)	$l_1$	$l_2$	$l_3$	...	$l_k$
$P_1$	$m/z_1$	$rt_1$	5%	$l_{11}$	0	$l_{13}$	...	$l_{1k}$
$P_2$	$m/z_2$	$rt_2$	8%	$l_{21}$	$l_{22}$	0	...	$l_{2k}$
$P_3$	$m/z_3$	$rt_3$	3%	0	$l_{32}$	$l_{33}$	...	$l_{3k}$
...	...	...	...	...	...	...	...	...
$P_n$	$m/z_n$	$rt_n$	2%	$l_{n1}$	$l_{n2}$	$l_{n3}$	...	$l_{nk}$

# Data analysis - Imputation

## Imputation with zeros

Alternative strategy: **Half-minimum imputation**

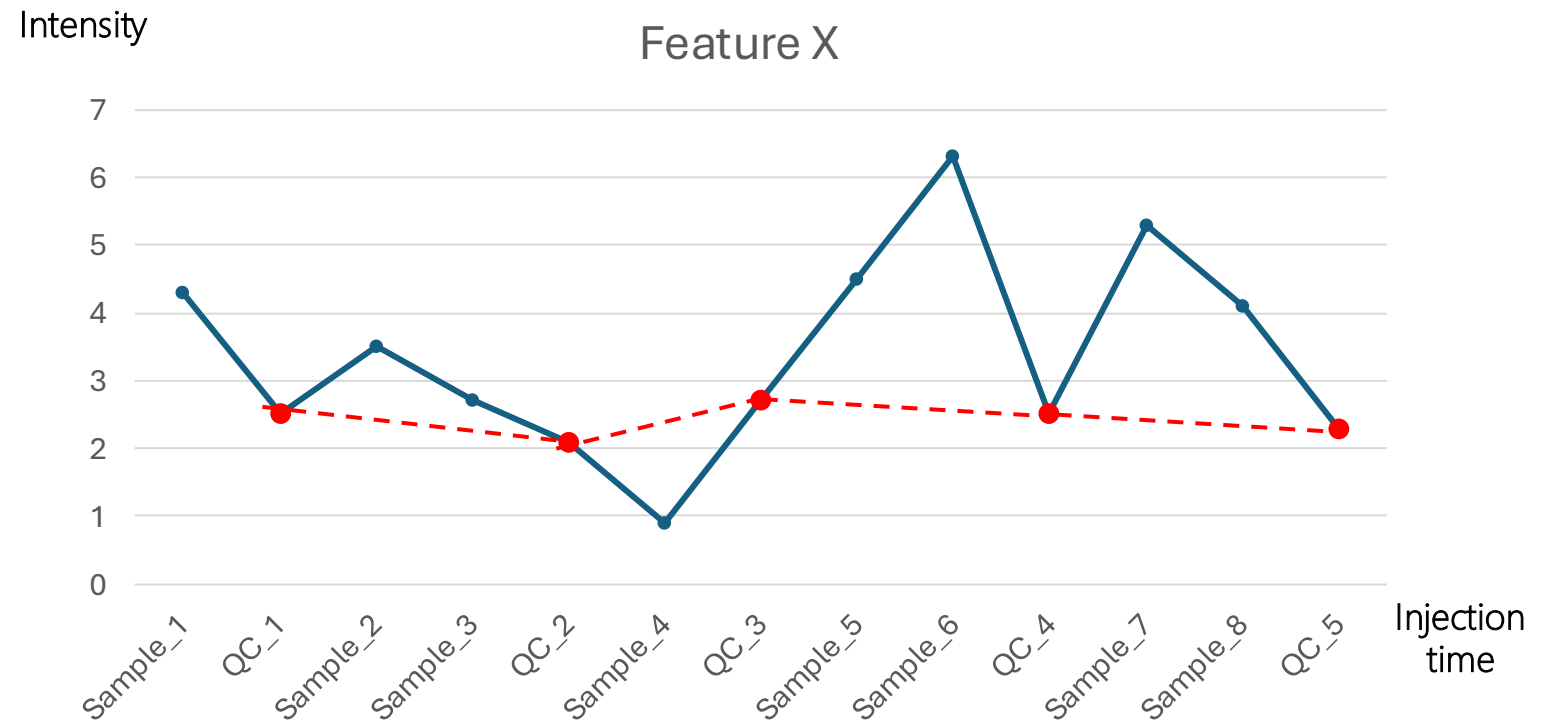
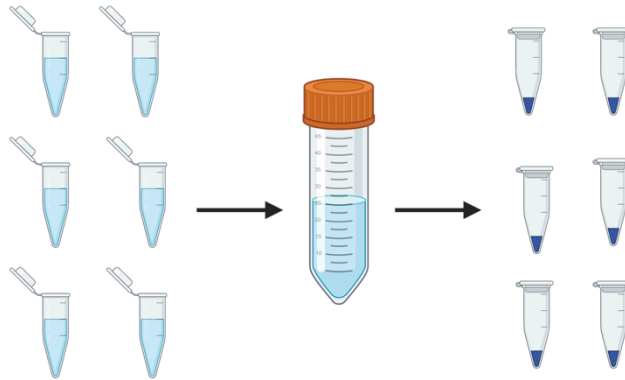
Assumption: Feature data is **Missing Not at Random (MNAR)** = metabolite could be just below the limit of detection  $\longrightarrow$  limit of detection is ...?  $\longleftarrow$

For each feature, find the minimum intensity value recorded, and divide it in half

Peak	m/z	rt	missingness (%)	$I_1$	$I_2$	$I_3$	...	$I_k$
$P_1$	$m/z_1$	$rt_1$	5%	$I_{11}$	$I_{22}/2$	$I_{13}$	...	$I_{1k}$
$P_2$	$m/z_2$	$rt_2$	8%	$I_{21}$	$I_{22}$	$I_{73}/2$	...	$I_{2k}$
$P_3$	$m/z_3$	$rt_3$	3%	$I_{51}/2$	$I_{32}$	$I_{33}$	...	$I_{3k}$
...	...	...	...	...	...	...	...	...
$P_n$	$m/z_n$	$rt_n$	2%	$I_{n1}$	$I_{n2}$	$I_{n3}$	...	$I_{nk}$

# Data analysis - Filtering and normalization

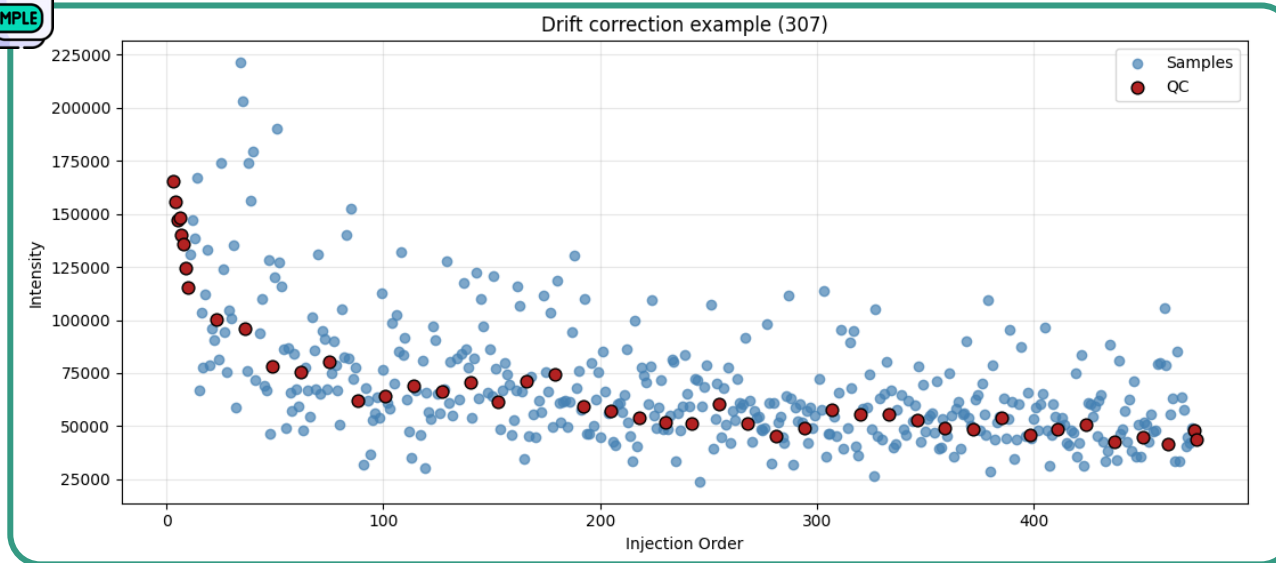
Drift correction



# Data analysis - Filtering and normalization

## Drift correction: LOESS

- We want to use QC samples to estimate artificial variation that comes from instrumental drift over time
- We want to subtract that variation from our data points

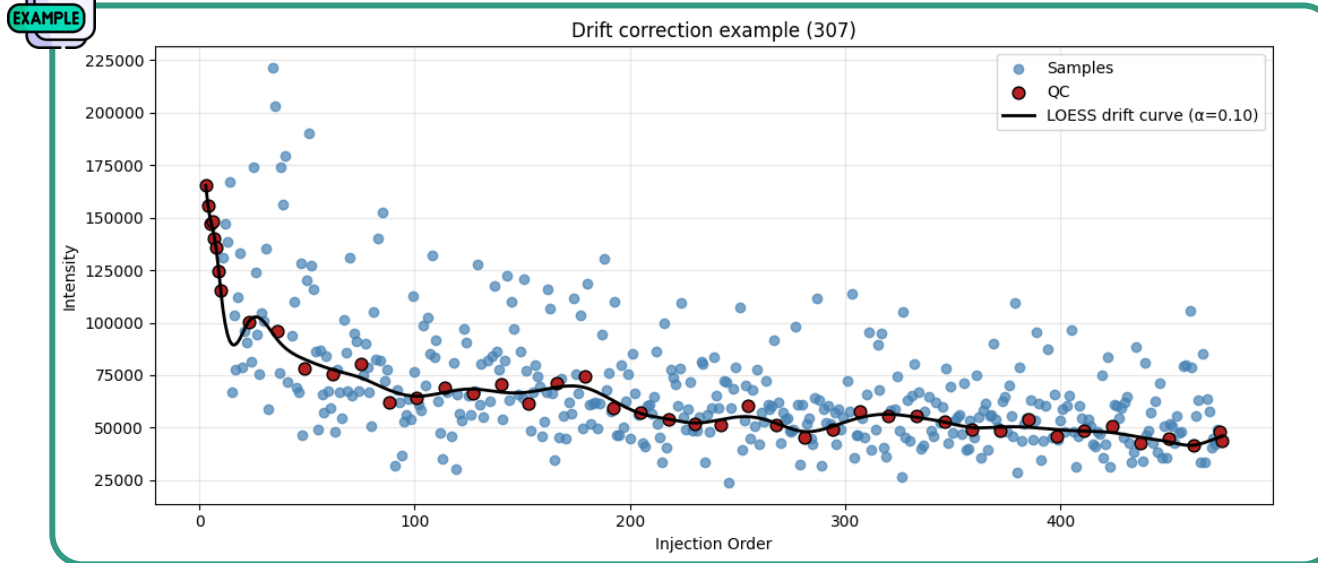


# Data analysis - Filtering and normalization

## Drift correction: LOESS

LOESS = Locally Estimated Scatterplot Smoothing

- Interpolate a value for QC samples at the run order value of every biological sample

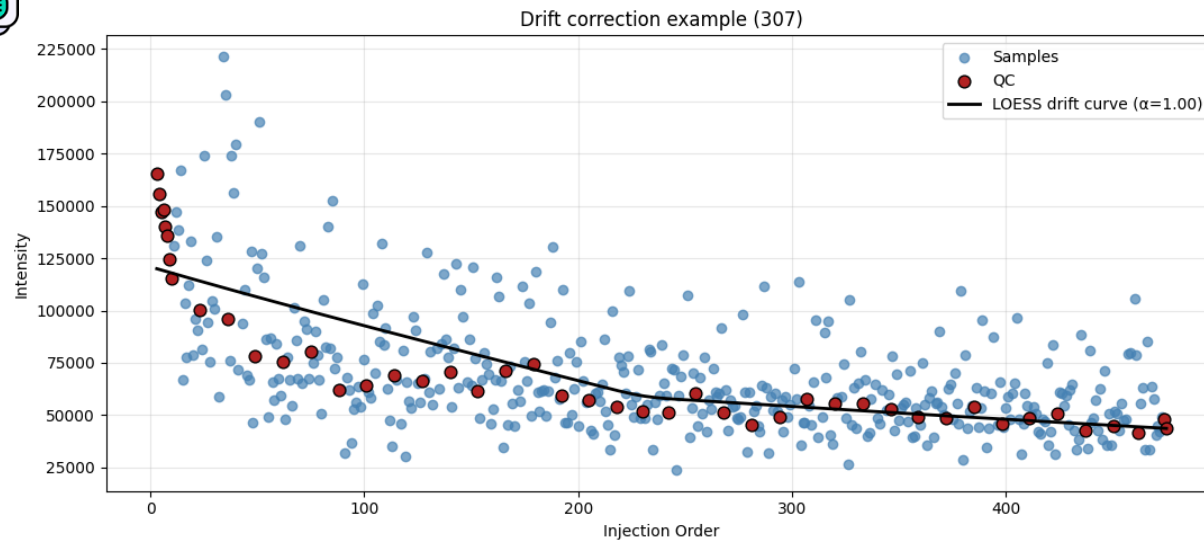


# Data analysis - Filtering and normalization

## Drift correction: LOESS

LOESS = Locally Estimated Scatterplot Smoothing

- Interpolate a value for QC samples at the run order value of every biological sample
- To reduce influence of extreme values, use LOESS Smoothing before interpolation

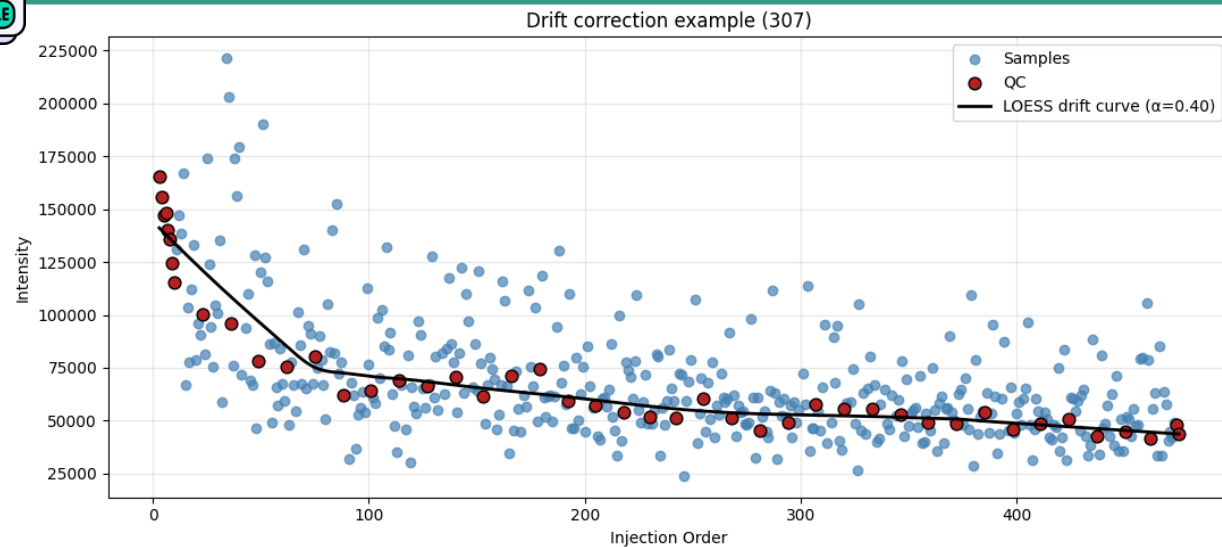


# Data analysis - Filtering and normalization

## Drift correction: LOESS

LOESS = Locally Estimated Scatterplot Smoothing

- Interpolate a value for QC samples at the run order value of every biological sample
- To reduce influence of extreme values, use LOESS Smoothing before interpolation
- Adjust level of smoothing with parameter alpha

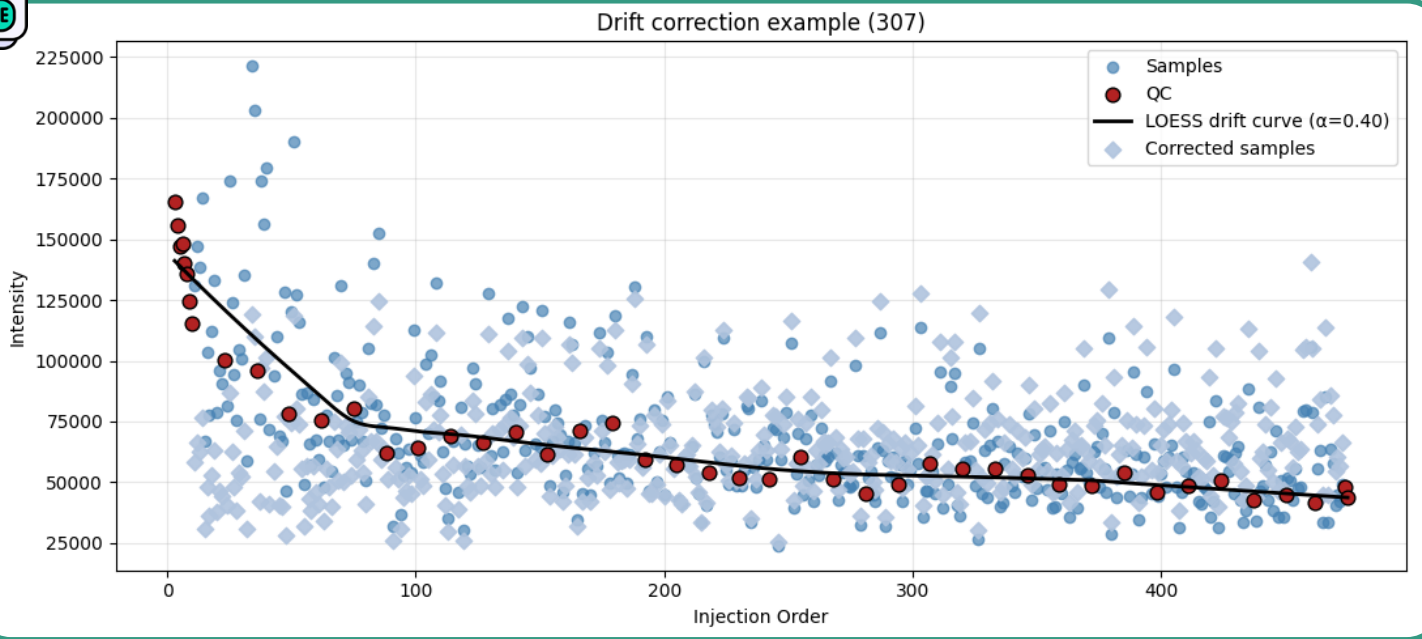


# Data analysis - Filtering and normalization

## Drift correction: LOESS

Now that we have an artificial QC value for each injection time point, can correct our biological sample values

$$\text{new values} = \frac{\text{original values}}{\text{drift curve}} \times \text{median(QCs)}$$



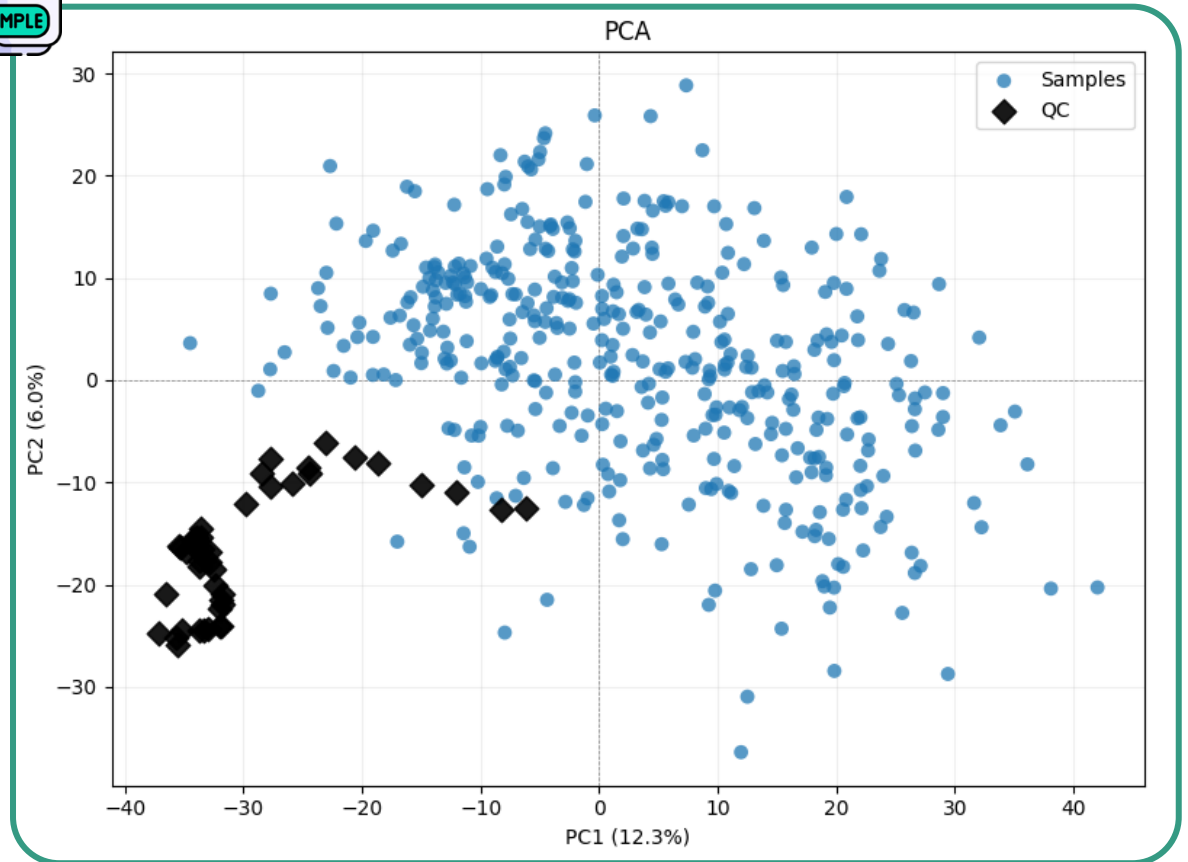
- Assumes drift is multiplicative: scales signal up or down
- Divide by drift curve to remove that factor
- Have relative intensity: how strong signal was relative to instrument at that moment
- Scale back to this specific feature's intensities with QC median

# Data analysis - Filtering and normalization

## Drift correction: CPCA

### Common Principal Components Analysis (CPCA)

- Finds variation that is common across all groups
- Assumption: technical variation (drift) is shared across all groups, whereas biological variation is unique
- Want to use CPCA to identify and then remove this variation



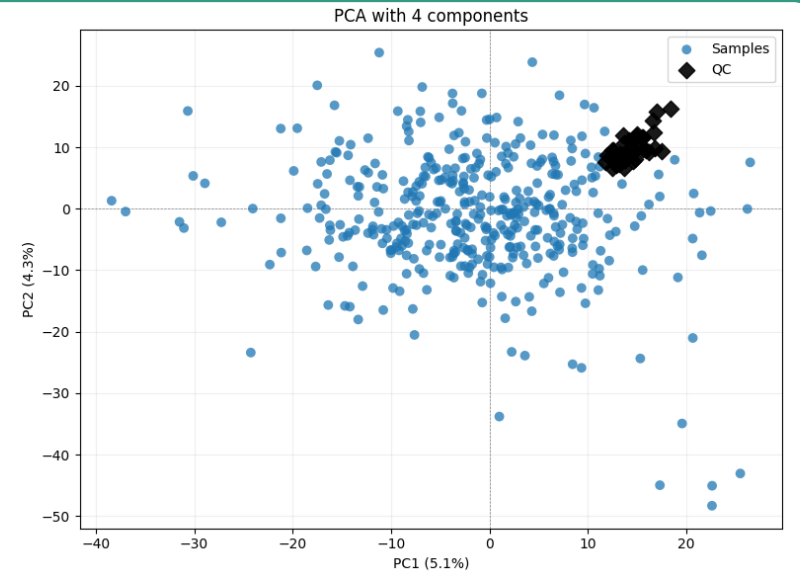
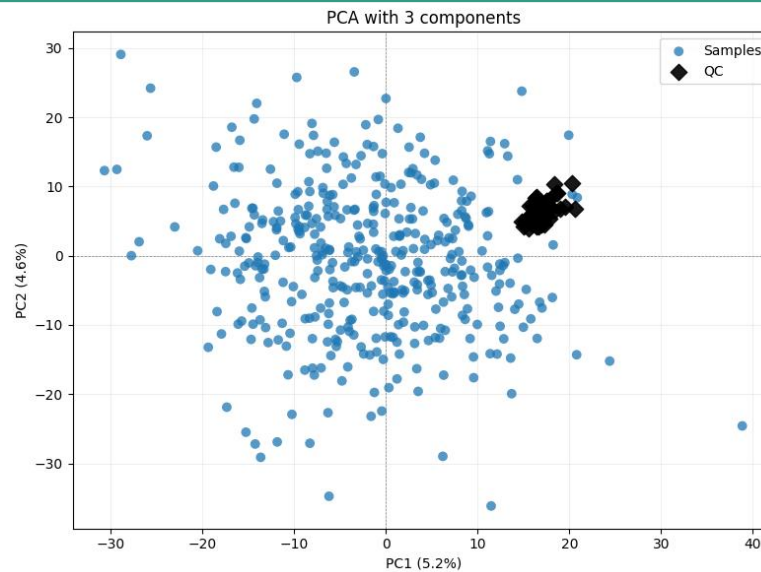
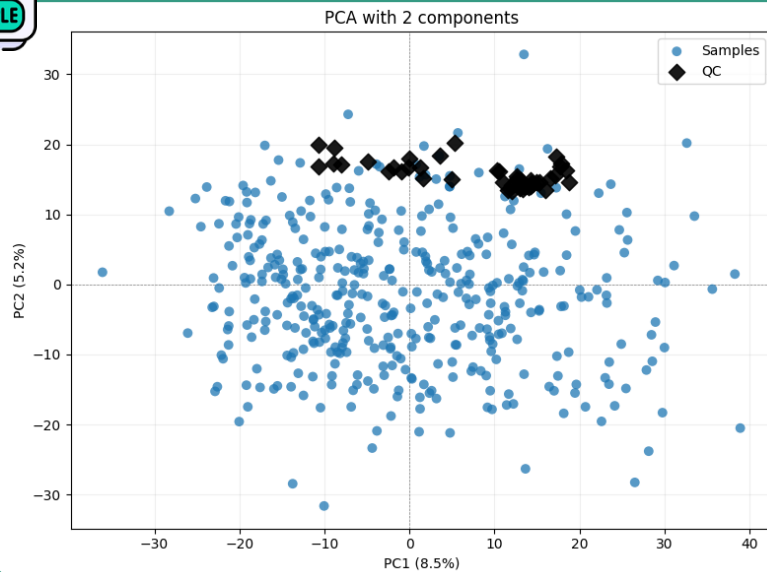
# Data analysis - Filtering and normalization

## Drift correction: CPCA

- Can try removing different numbers of components
- Choose the one where the cluster of QCs after removal of principal component is the tightest
- Can calculate the distance of points to cluster centroid to assess objectively



EXAMPLE

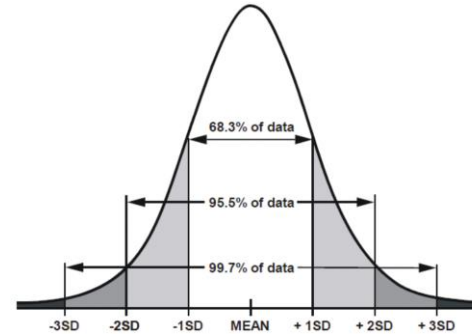


# Data analysis - Filtering and normalization

## Normalization

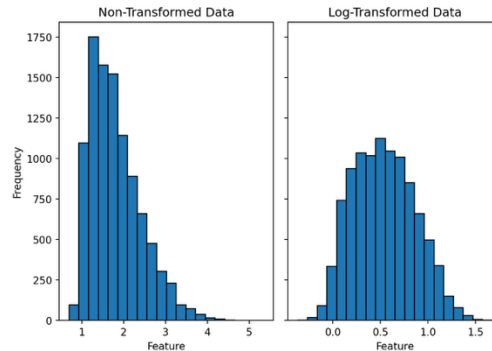
Z-score normalization  
(scaling)

$$z = \frac{x - \mu}{\sigma}$$



- Metabolite intensities can have vastly different magnitudes
- Gives every feature equal weight

Log transformation



- Metabolomics data is right-skewed (lower bounded at 0, but no upper bound)
- Stabilizes variance
- Makes fold changes more interpretable

Internal standards

$$\text{Normalized intensity} = \frac{\text{Analyte intensity}}{\text{Internal standard intensity}}$$

- Known amount
- Usually not used in untargeted metabolomics, there are too many metabolites

# Data analysis - Statistical Analysis

## Univariate analyses – one metabolite at a time

### Comparing two groups

- T-test on log-transformed data (parametric)
- Mann-Whitney U test, if normality can't be assumed (non-parametric)

### Comparing more than two groups

- ANOVA (ANalysis Of VAriance) (parametric)
- Kruskal-Wallis, if normality can't be assumed (non-parametric)

### Continuous outcomes

- Linear regression
- Correlation

### Adjusting for covariates

- ANCOVA (ANalysis of COVAriance)
- Multiple linear regression

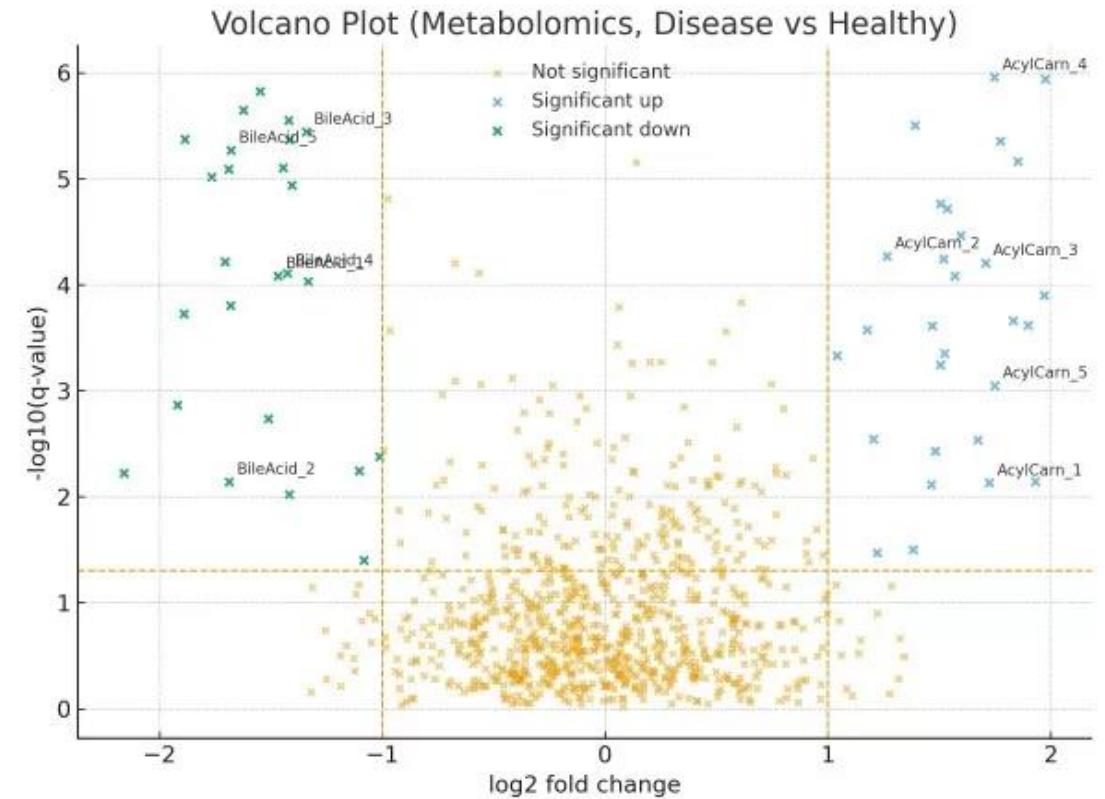
# Data analysis - Statistical Analysis

Univariate analyses – one metabolite at a time

Fold changes = difference in the average metabolite level between sample groups

Important: multiple testing correction

- Huge numbers of features, so p-values need to be corrected
- Common: Benjamini-Hochberg (BH) False Discovery Rate
- Rather consider adjusted p-values (or q-values)



# Data analysis - Statistical Analysis

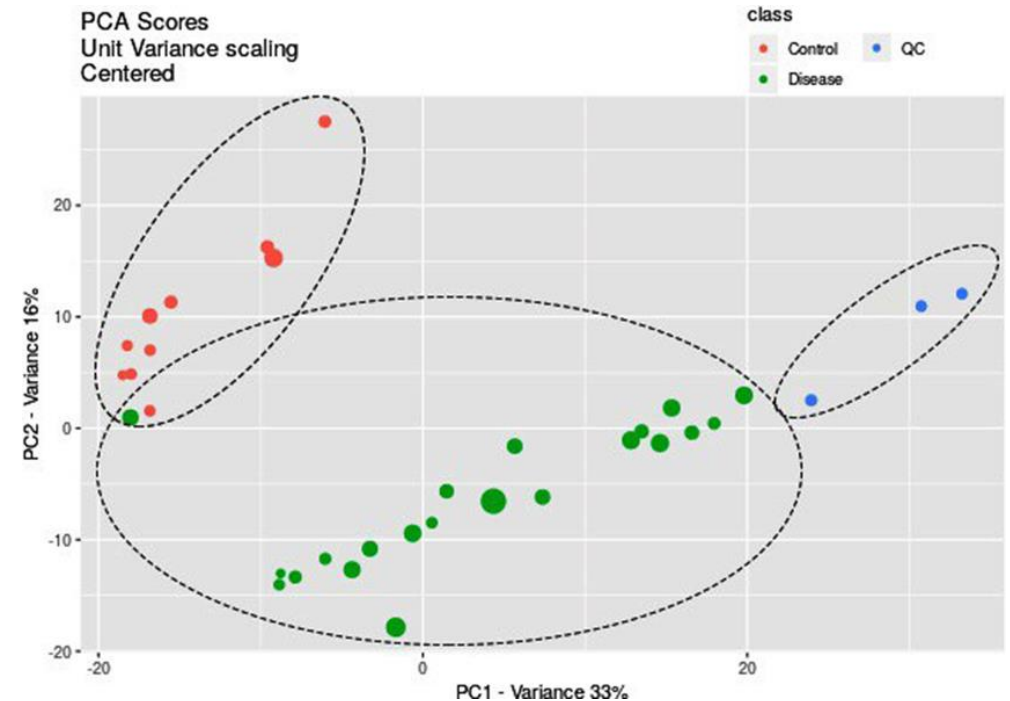
Multivariate analysis – all metabolites together

## Principal Component Analysis (PCA)

- Reduces high dimensionality of metabolomics data into principal components (PCs)
- PCs capture largest sources of variance
- Can use this to see if controls cluster well

## Partial Least-Squares Discriminant Analysis (PLS-DA)

- Discovers which metabolites drive the most separation between groups
- Latent variables are rotated specifically to discriminate between classes



# Data analysis – Example Pipeline

## 1. Filtering

- Modified 80%-rule
- Blanks filtering
- CV-filtering

## 2. Imputation

- Half-minimum imputation

## 3. Drift correction

- LOESS smoothing-based correction & CPCA
- Assess which one works better

& using PCAs throughout to check data

## 4. Normalization

- Log transforming data

## 5. Statistical analysis

- ANOVA
- Volcano plot



**Break!**

# Data analysis – Hands on

# Metabolomics Interpretation

## Peak annotation vs peak identification

Identification levels according to the Chemical Analysis Working Group of the **Metabolomics Standards Initiative**:

Level 1 – **Identified** compounds require comparison of at least two independent properties (e.g., retention time and MS/MS spectra) with those of an authentic standard analysed under identical conditions.

# Metabolomics Interpretation

## Peak annotation vs peak identification

Identification levels according to the Chemical Analysis Working Group of the **Metabolomics Standards Initiative**:

Level 1 – **Identified** compounds require comparison of at least two independent properties (e.g., retention time and MS/MS spectra) with those of an authentic standard analysed under identical conditions.

Level 2 – Putatively **annotated** compounds based on similarity to MS/MS data in commercial or public databases (reference standards are not available).

# Metabolomics Interpretation

## Peak annotation vs peak identification

Identification levels according to the Chemical Analysis Working Group of the **Metabolomics Standards Initiative**:

Level 1 – **Identified** compounds require comparison of at least two independent properties (e.g., retention time and MS/MS spectra) with those of an authentic standard analysed under identical conditions.

Level 2 – Putatively **annotated** compounds based on similarity to MS/MS data in commercial or public databases (reference standards are not available).

Level 3 – Putatively characterized **compound class**.

Level 4 – **Unknown** feature.

# Metabolomics Interpretation

## Peak annotation vs peak identification

Identification levels according to the Chemical Analysis Working Group of the **Metabolomics Standards Initiative**:

Level 1 – **Identified** compounds require comparison of at least two independent properties (e.g., retention time and MS/MS spectra) with those of an authentic standard analysed under identical conditions.

Level 2 – Putatively **annotated** compounds based on similarity to MS/MS data in commercial or public databases (reference standards are not available).

Level 3 – Putatively characterized **compound class**.

Level 4 – **Unknown** feature.

LC-MS<sup>1</sup>

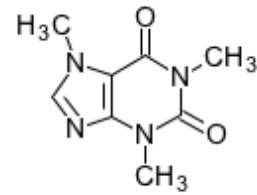
# Metabolomics Interpretation

## Peak annotation using LC-MS<sup>1</sup>

Peak annotation is a **major bottleneck** and direct matching of masses is **not enough**:

1

A single metabolite produces multiple peaks due to adducts and fragments.



C07481

Caffeine

Exact mass = 194.0804

— LC-MS + —>

[M+H]<sup>+</sup>

[M+Na]<sup>+</sup>

[M+K]<sup>+</sup>

[M+NH<sub>4</sub>]<sup>+</sup>

# Metabolomics Interpretation

Peak annotation using LC-MS<sup>1</sup>

Peak annotation is a **major bottleneck** and direct matching of masses is **not enough**:

2

DB matching results in hundreds of potential candidates.

m/z = 195.0877



Positive mode  
10 different adducts

635 Results:

Caffeine [M+H]<sup>+</sup> - mass = 194.0804

Enprofylline [M+H]<sup>+</sup> - mass = 194.0804

Luminol [M+NH<sub>4</sub>]<sup>+</sup> - mass = 177.0538

...

# Metabolomics Interpretation

## Peak annotation using LC-MS<sup>1</sup>

Peak annotation is a **major bottleneck** and direct matching of masses is **not enough**:

1

A single metabolite produces multiple peaks due to adducts and fragments.

2

DB matching results in hundreds of potential candidates.

3

The metabolome is not completely annotated.



Computational algorithms that use m/z, rt and intensity to return putative annotations.

# Metabolomics Interpretation

Peak annotation using LC-MS<sup>1</sup>

Peak annotation is a **major bottleneck** and direct matching of masses is **not enough**:

1

A single metabolite produces multiple peaks due to adducts and fragments.

2

DB matching results in hundreds of potential candidates.

3

The metabolome is not completely annotated.

Computational algorithms



Low confidence level

putative annotations.

# Metabolomics Interpretation

## Peak annotation vs peak identification

Identification levels according to the Chemical Analysis Working Group of the Metabolomics Standards Initiative:

Level 1 – **Identified** compounds require comparison of at least two independent properties (e.g., retention time and MS/MS spectra) with those of an authentic standard analysed under identical conditions.

Level 2 – Putatively **annotated** compounds based on similarity to MS/MS data in commercial or public databases (reference standards are not available).

LC-MS<sup>2</sup>

Level 3 – Putatively characterized **compound class**.

Level 4 – **Unknown** feature.

# Metabolomics Interpretation

Peak annotation using LC-MS<sup>2</sup> spectra

1. Spectral library matching
2. In-silico fragmentation
3. Fragmentation tree methods

# Metabolomics Interpretation

## Peak annotation using LC-MS<sup>2</sup> spectra

### 1. Spectral library matching

- Match our query spectrum with known MS<sup>2</sup> spectra



We only have annotated spectra for a fraction of the known metabolome

We don't know the whole metabolome

### 2. In-silico fragmentation

### 3. Fragmentation tree methods

# Metabolomics Interpretation

## Peak annotation using LC-MS<sup>2</sup> spectra

### 1. Spectral library matching

- Match our query spectrum with known MS<sup>2</sup> spectra



We only have annotated spectra for a fraction of the known metabolome

We don't know the whole metabolome

### 2. In-silico fragmentation

- Predicting MS<sup>2</sup> spectra from chemical structures to increase spectral libraries



Still very unreliable, especially for large complex molecules and lipids

### 3. Fragmentation tree methods

# Metabolomics Interpretation

## Peak annotation using LC-MS<sup>2</sup> spectra

### 1. Spectral library matching

- Match our query spectrum with known MS<sup>2</sup> spectra



We only have annotated spectra for a fraction of the known metabolome

We don't know the whole metabolome

### 2. In-silico fragmentation

- Predicting MS<sup>2</sup> spectra from chemical structures to increase spectral libraries



Still very unreliable, especially for large complex molecules and lipids

### 3. Fragmentation tree methods

- To predict molecular formulas, fingerprints or even novel structures.



Excellent results for smaller metabolites (<400 Da), but unreliable for larger molecules

# Metabolomics Interpretation

## Peak annotation using LC-MS<sup>2</sup> spectra

### 1. Spectral library matching

- Match our query spectrum with known MS<sup>2</sup> spectra



We only have annotated spectra for a fraction of the known metabolome

We don't know the whole metabolome

### 2. In-silico fragmentation

- Predicting MS<sup>2</sup> spectra from chemical structures to increase spectral libraries



Still very unreliable, especially for large complex molecules and lipids

### 3. Fragmentation tree methods

- To predict molecular formulas, fingerprints or even novel structures.

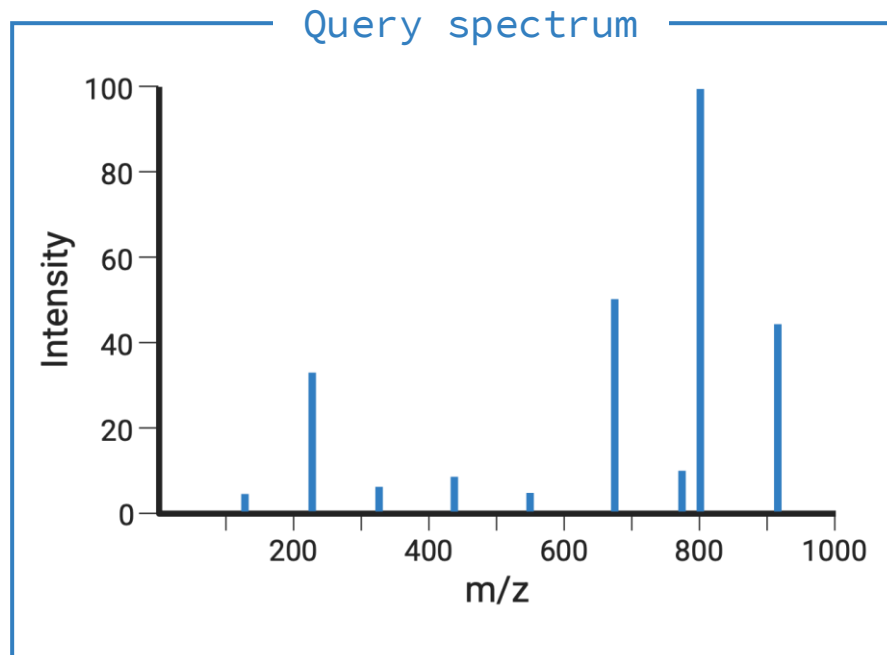


Excellent results for smaller metabolites (<400 Da), but unreliable for larger molecules

# Metabolomics Interpretation

Peak annotation using LC-MS<sup>2</sup> spectra

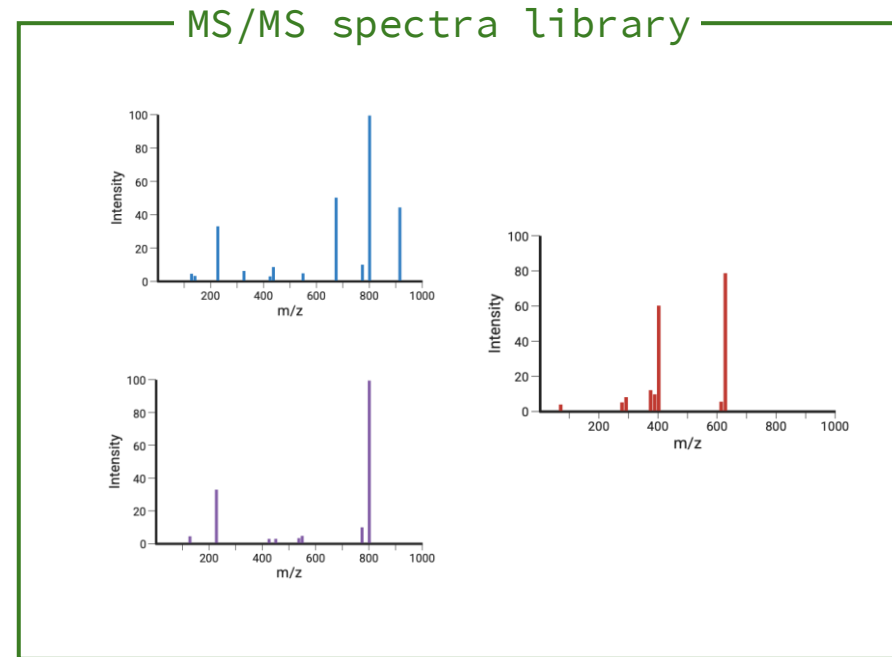
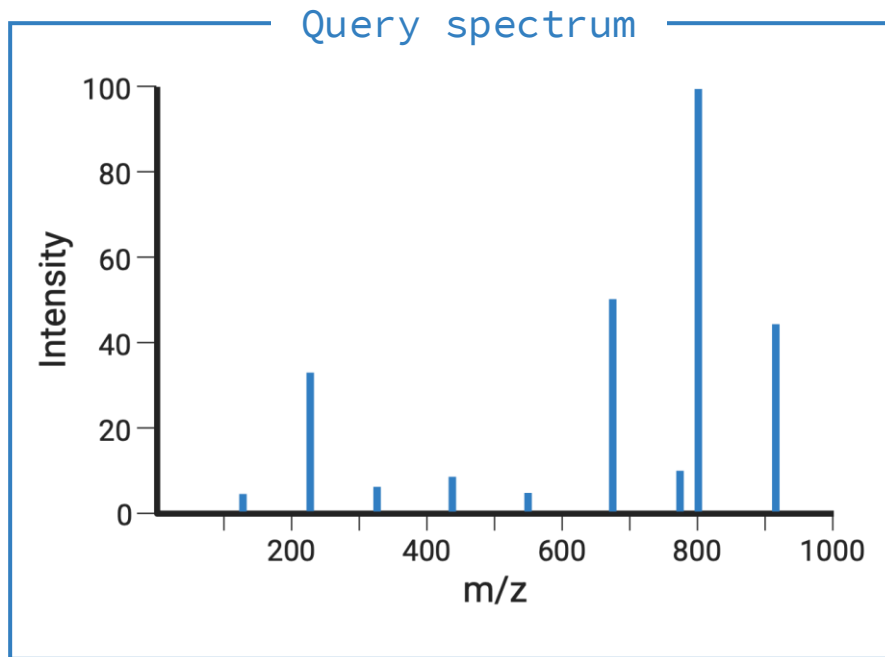
Spectral library matching



# Metabolomics Interpretation

Peak annotation using LC-MS<sup>2</sup> spectra

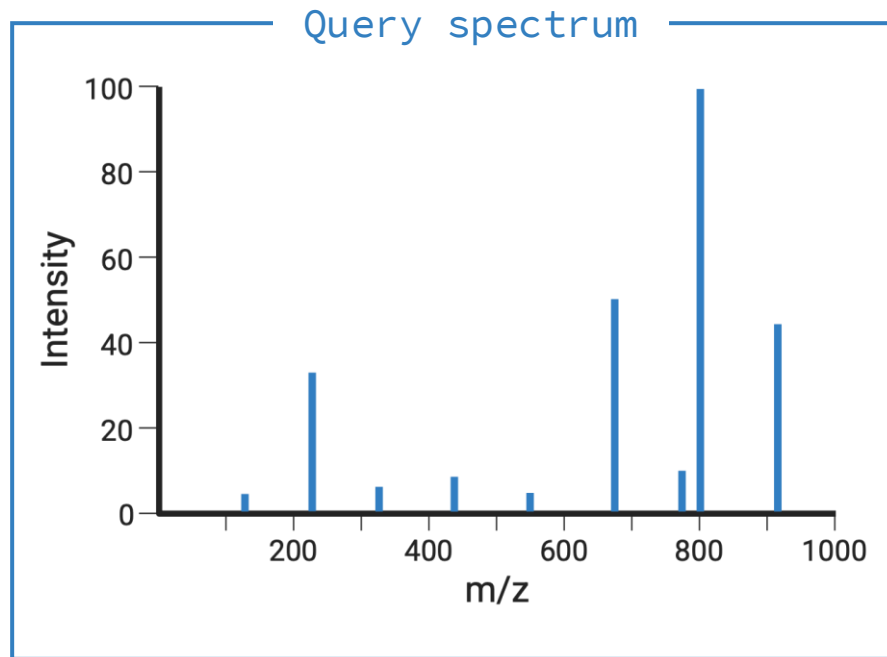
## Spectral library matching



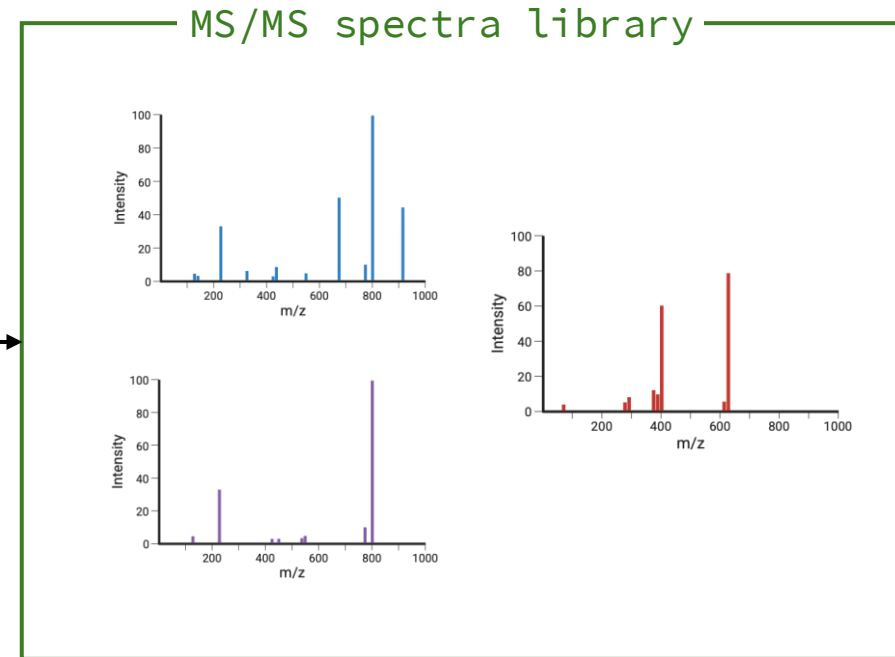
# Metabolomics Interpretation

Peak annotation using LC-MS<sup>2</sup> spectra

## Spectral library matching



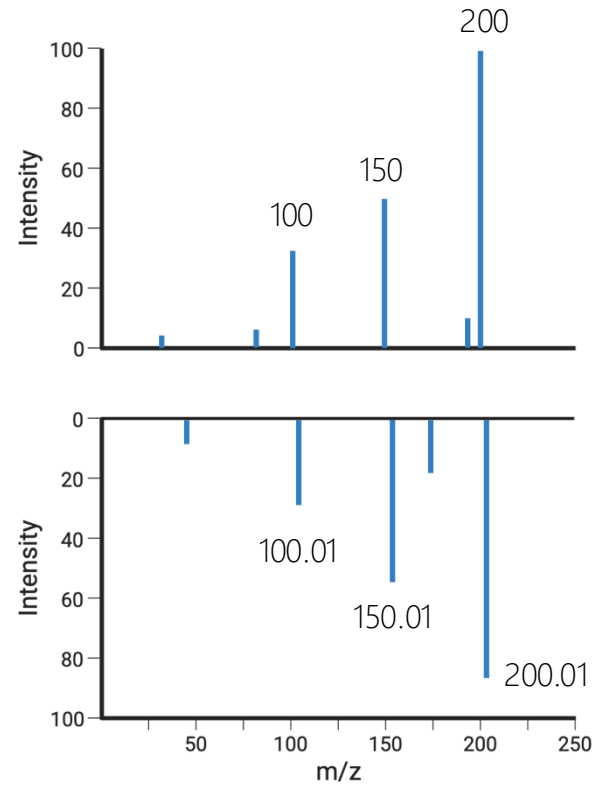
← Similarity →



# Metabolomics Interpretation

Peak annotation using LC-MS<sup>2</sup> spectra

## Cosine similarity



Query spectrum

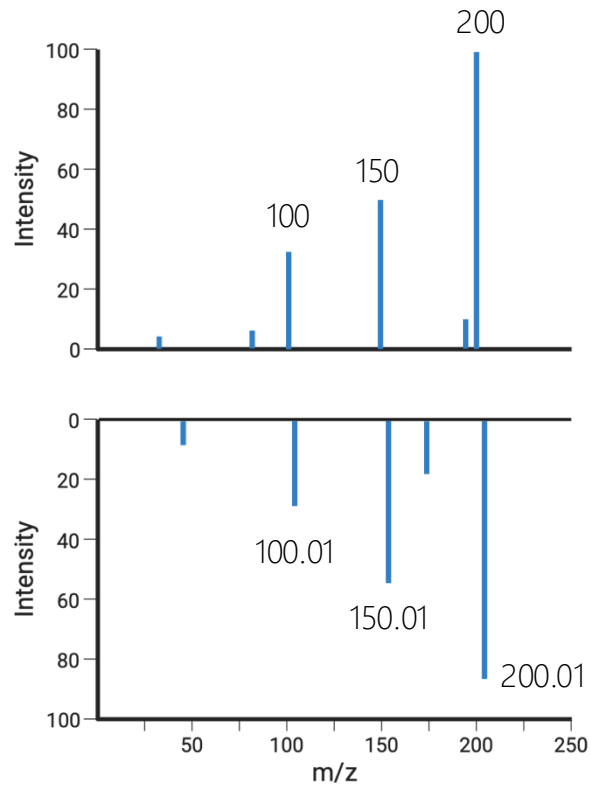
Reference spectrum

# Metabolomics Interpretation

Peak annotation using LC-MS<sup>2</sup> spectra

## Cosine similarity

1. Find peaks that match:  $|mz(p) - mz(p')| < t$  (e.g., 0.05 Da)

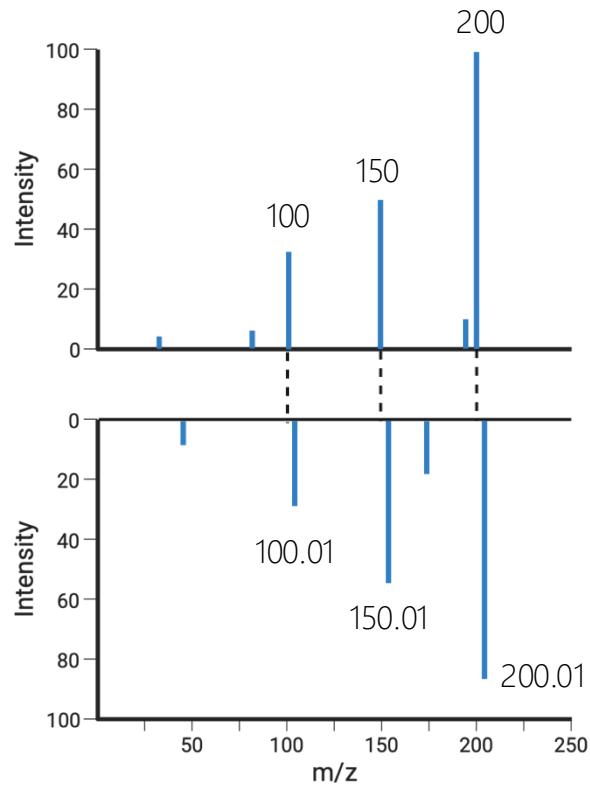


# Metabolomics Interpretation

Peak annotation using LC-MS<sup>2</sup> spectra

## Cosine similarity

1. Find peaks that match:  $|mz(p) - mz(p')| < t$  (e.g., 0.05 Da)



# Metabolomics Interpretation

Peak annotation using LC-MS<sup>2</sup> spectra

## Cosine similarity

1. Find peaks that match:  $|mz(p) - mz(p')| < t$  (e. g., 0.05 Da)
2. Intensities of matching peaks are converted to vectors

$I = [38, 46, 101]$

$I' = [30, 59, 85]$

# Metabolomics Interpretation

Peak annotation using LC-MS<sup>2</sup> spectra

## Cosine similarity

1. Find peaks that match:  $|mz(p) - mz(p')| < t$  (e. g., 0.05 Da)

$I = [38, 46, 101]$

2. Intensities of matching peaks are converted to vectors

3. Cosine similarity:

$I' = [30, 59, 85]$

$$\cos(\theta) = \frac{A \cdot B}{\|A\| \|B\|} = \frac{\sum_{i=1}^n A_i B_i}{\sqrt{\sum_{i=1}^n A_i^2} \sqrt{\sum_{i=1}^n B_i^2}} = 0.98$$

# Metabolomics Interpretation

Peak annotation using LC-MS<sup>2</sup> spectra

What if ... Two molecules differ by a methyl group:

Molecule A precursor = 300 Da

Molecule B precursor = 314 Da



Difference = 14 Da

Many fragments may shift +14 Da

# Metabolomics Interpretation

Peak annotation using LC-MS<sup>2</sup> spectra

What if ... Two molecules differ by a methyl group:

Molecule A precursor = 300 Da

Molecule B precursor = 314 Da



Difference = 14 Da

Many fragments may shift +14 Da

— Cosine similarity  
approach →

No matching peaks

Similarity = 0

# Metabolomics Interpretation

Peak annotation using LC-MS<sup>2</sup> spectra

What if ... Two molecules differ by a methyl group:

Molecule A precursor = 300 Da

Molecule B precursor = 314 Da



Difference = 14 Da

Many fragments may shift +14 Da

— Cosine similarity approach →

No matching peaks

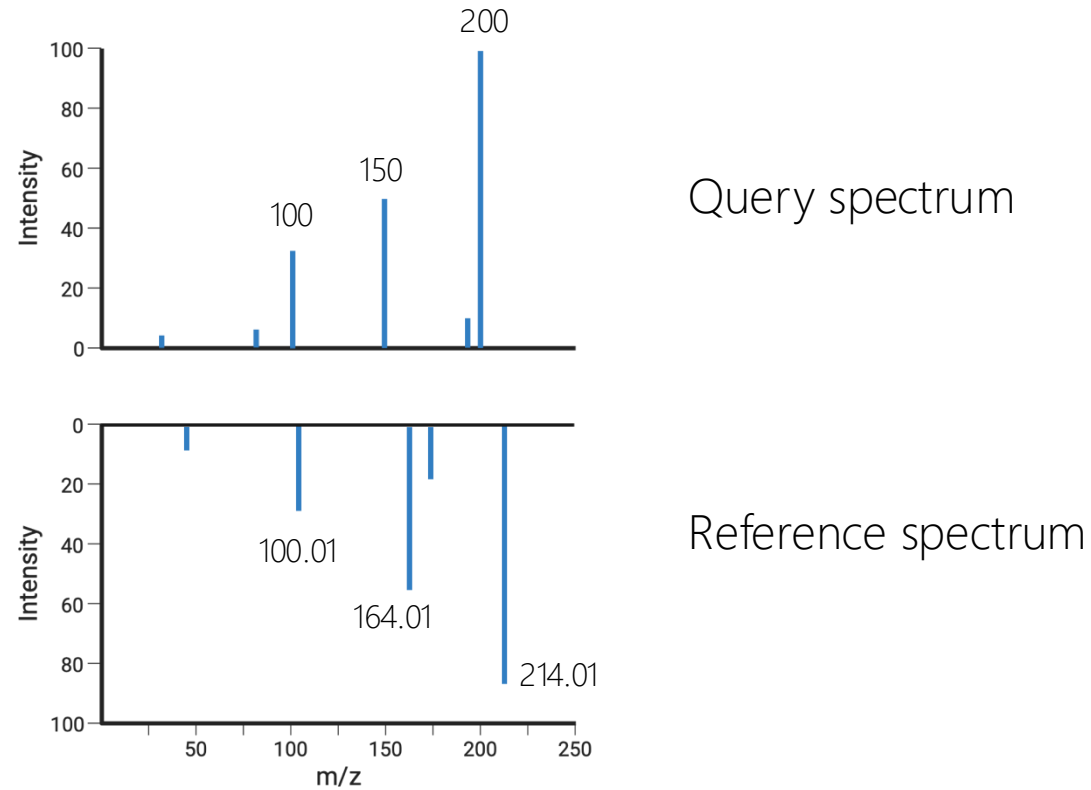
Similarity = 0

But molecules are structurally very similar... Modified cosine similarity

# Metabolomics Interpretation

Peak annotation using LC-MS<sup>2</sup> spectra

## Modified cosine similarity



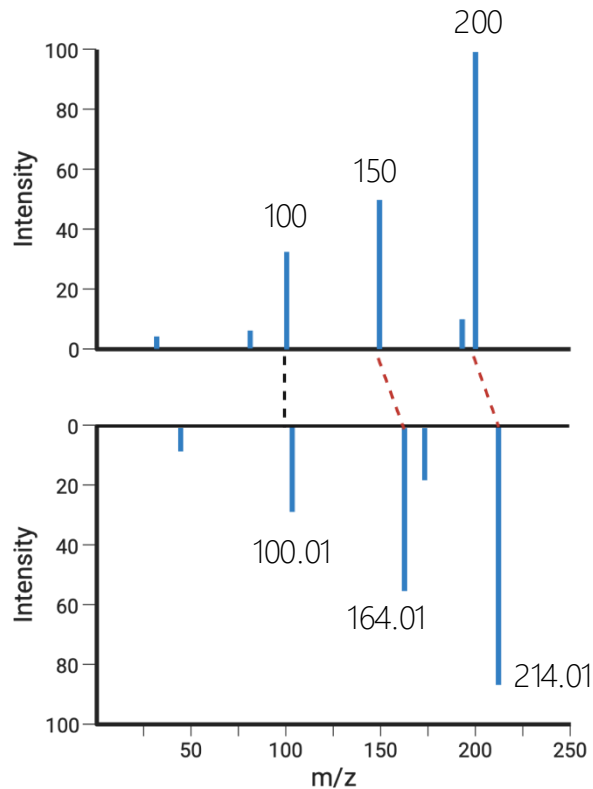
# Metabolomics Interpretation

Peak annotation using LC-MS<sup>2</sup> spectra

## Modified cosine similarity

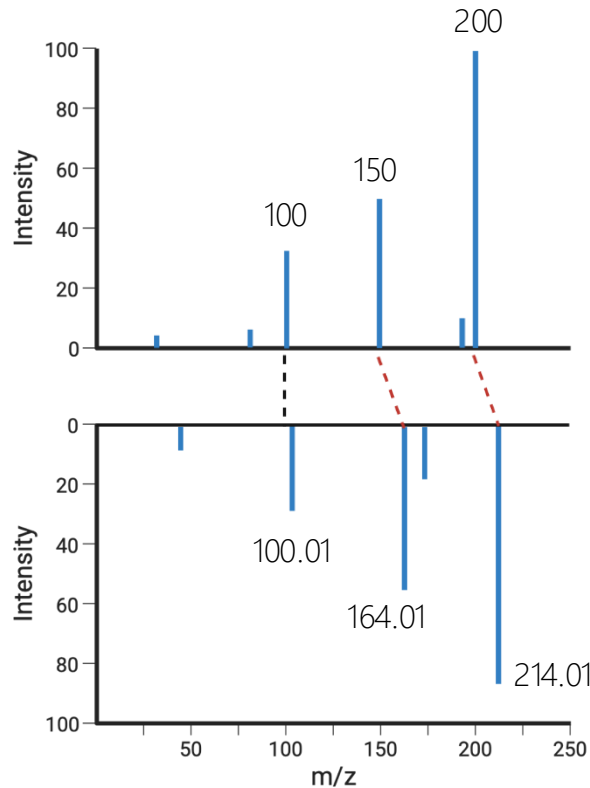
1. Find peaks that match:  $|mz(p) + M - mz(p')| < t$

Precursor mass difference  $M = PM(S') - PM(S) = 14$  Da



# Metabolomics Interpretation

Peak annotation using LC-MS<sup>2</sup> spectra



## Modified cosine similarity

1. Find peaks that match:  $|mz(p) + M - mz(p')| < t$   
Precursor mass difference  $M = PM(S') - PM(S) = 14$  Da
2. Intensities of matching peaks are converted to vectors
3. Cosine similarity:

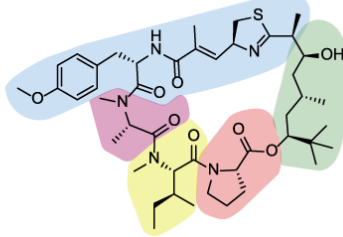
$$\cos(\theta) = \frac{A \cdot B}{\|A\| \|B\|} = \frac{\sum_{i=1}^n A_i B_i}{\sqrt{\sum_{i=1}^n A_i^2} \sqrt{\sum_{i=1}^n B_i^2}}$$

# Metabolomics Interpretation

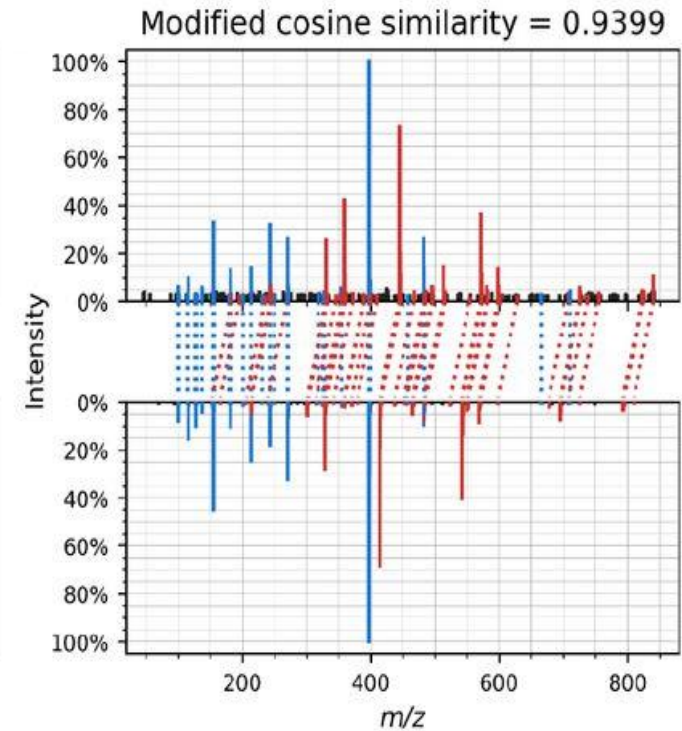
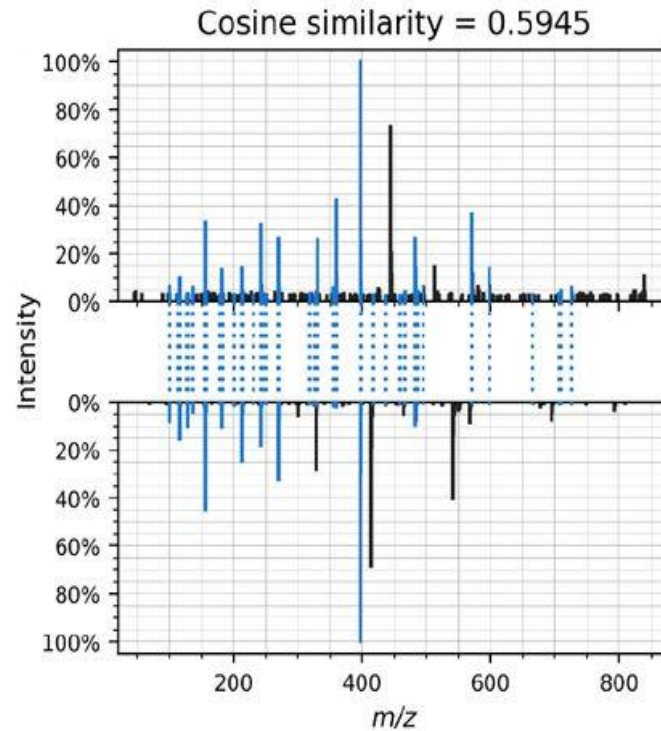
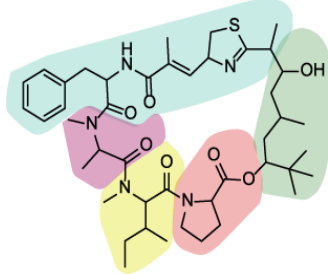
Peak annotation using LC-MS<sup>2</sup> spectra

## Spectral library matching

▲ Apratoxin A ( $m/z = 840.497$ )



◆ Apratoxin A - 30.010 Da ( $m/z = 810.487$ )



# Metabolomics Interpretation

Peak annotation using LC-MS<sup>2</sup> spectra

Spectral library matching

However...

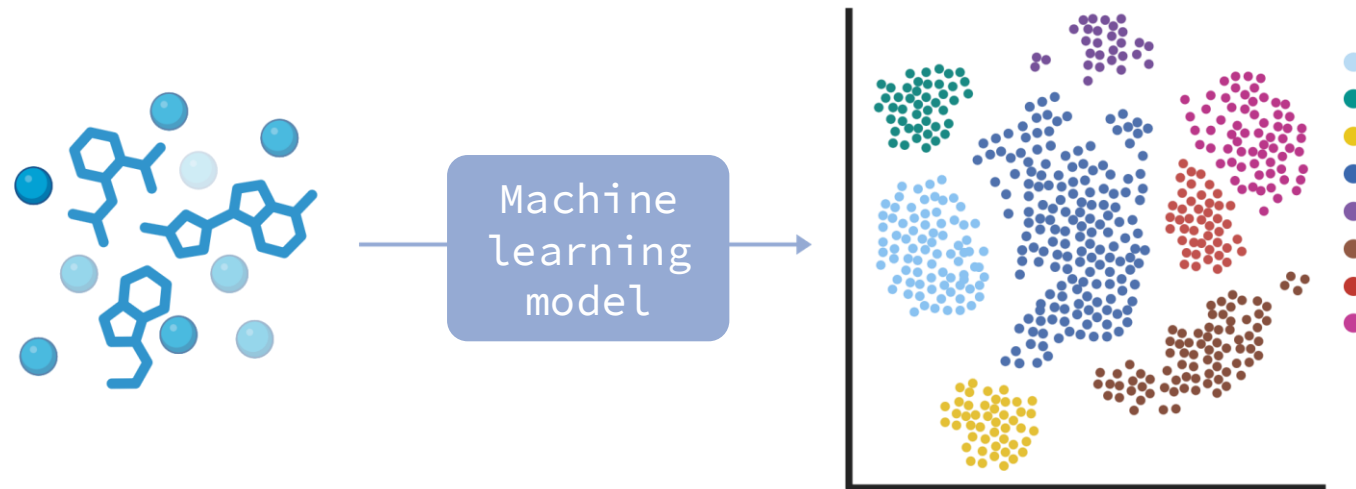
Small chemical modifications → **Low** modified cosine scores

Ideally, we want an approach that understands chemical structure!

# Metabolomics Interpretation

Peak annotation using LC-MS<sup>2</sup> spectra

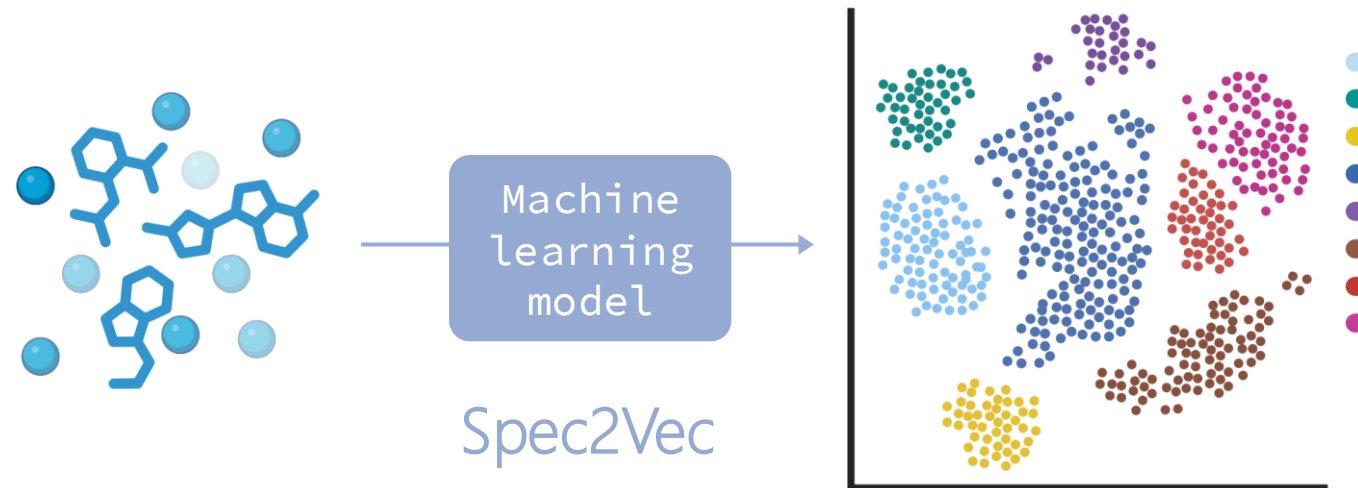
Machine learning-based similarities



# Metabolomics Interpretation

Peak annotation using LC-MS<sup>2</sup> spectra

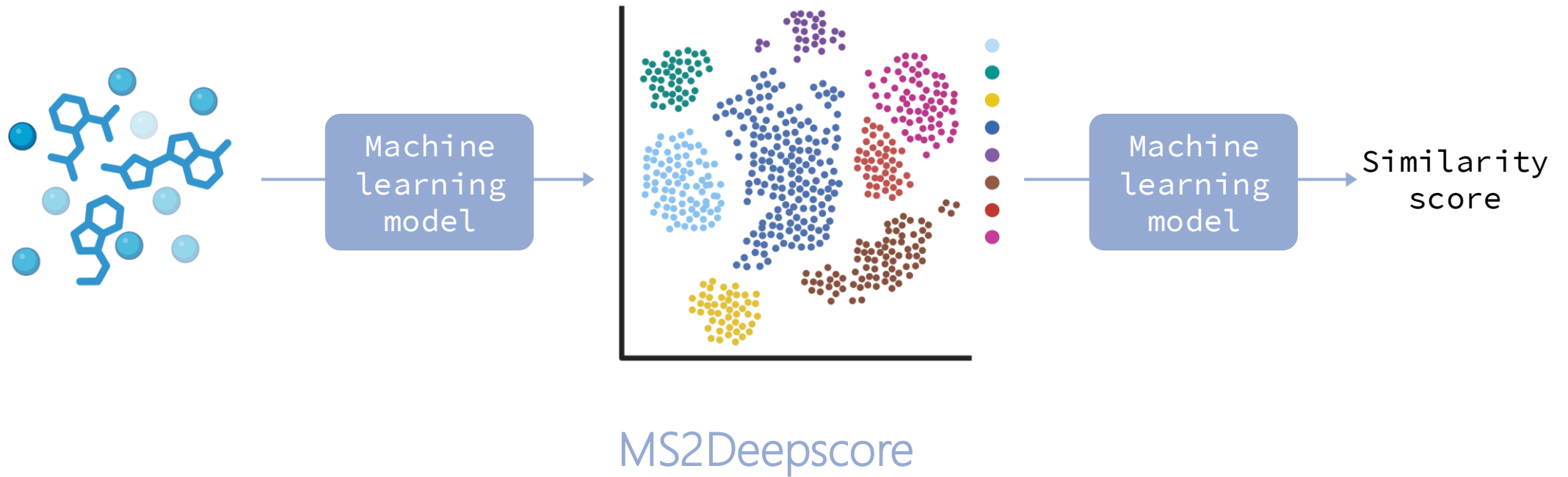
Machine learning-based similarities



# Metabolomics Interpretation

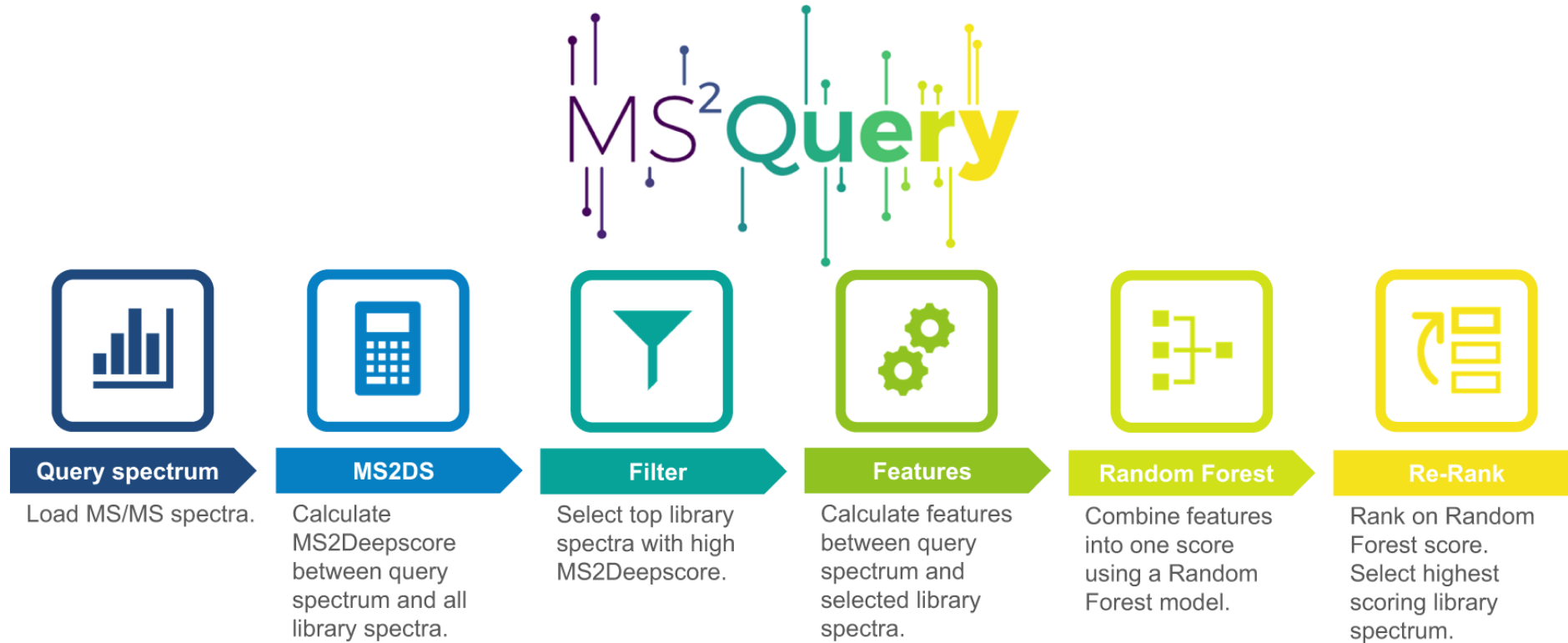
Peak annotation using LC-MS<sup>2</sup> spectra

Machine learning-based similarities



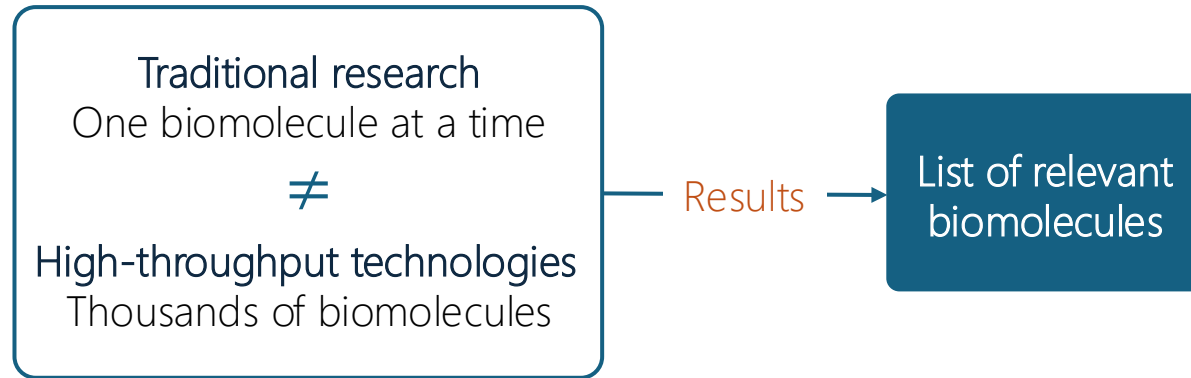
# Metabolomics Interpretation

Peak annotation with MS2Query



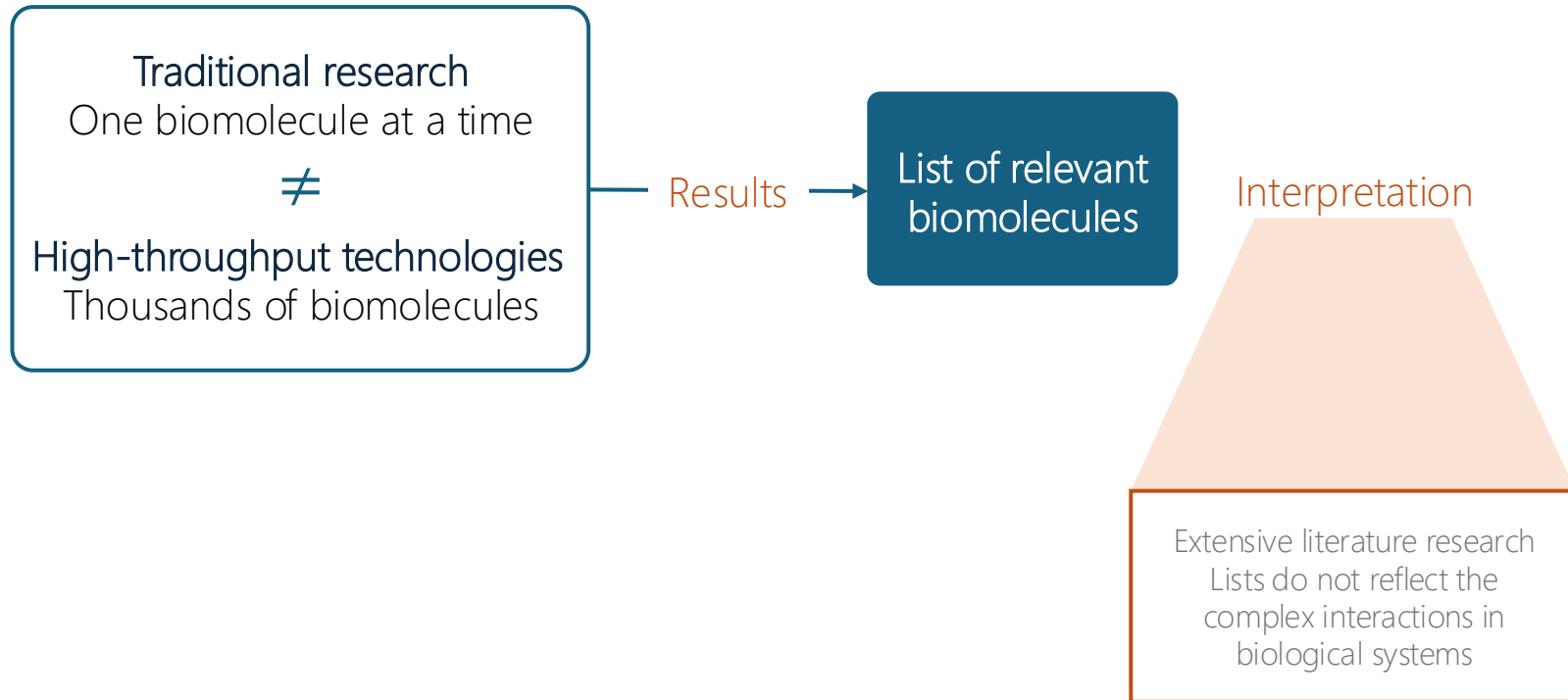
# Metabolomics Interpretation

Pathway enrichment



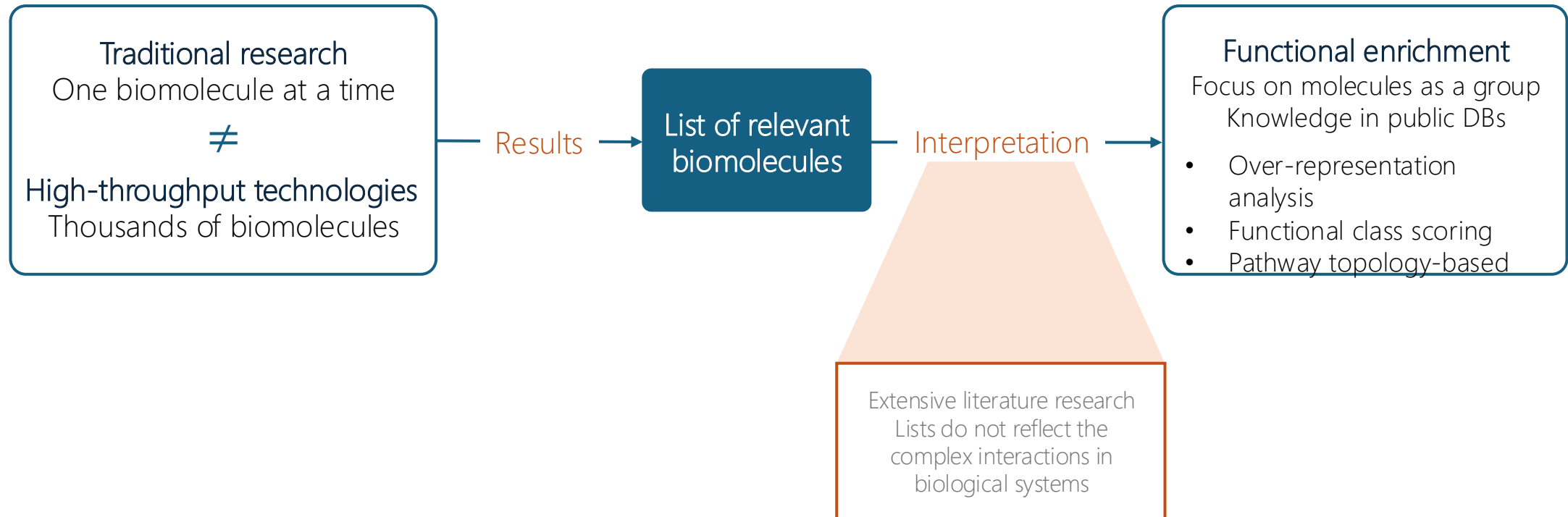
# Metabolomics Interpretation

Pathway enrichment



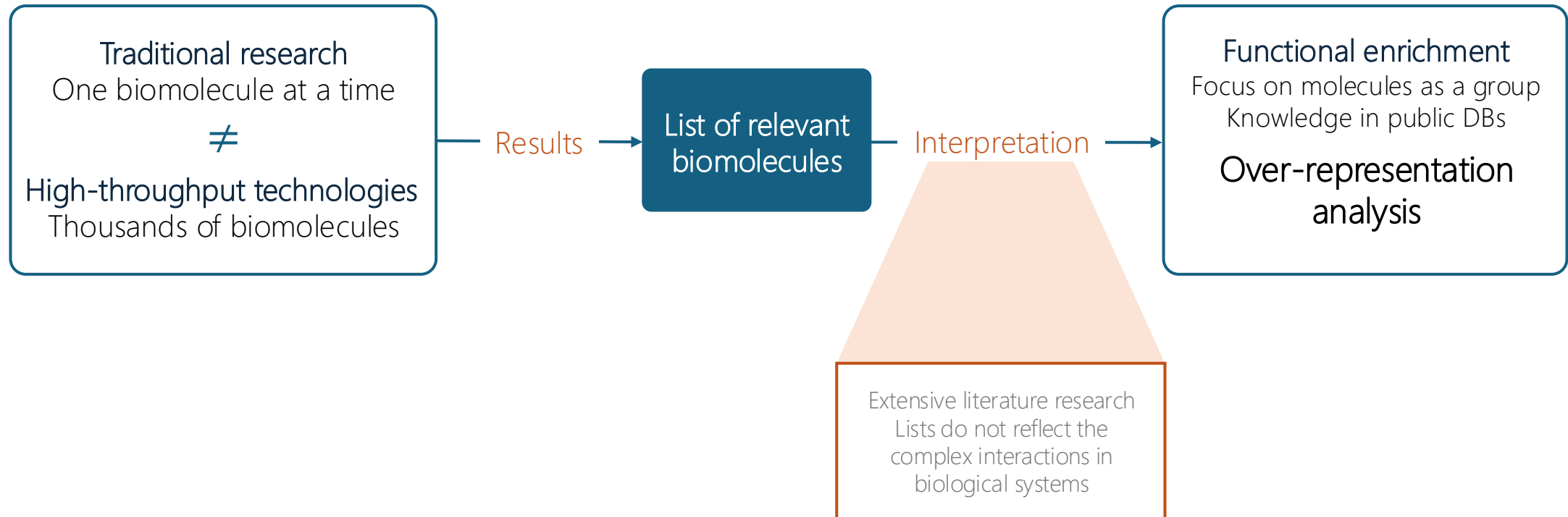
# Metabolomics Interpretation

## Pathway enrichment



# Metabolomics Interpretation

## Pathway enrichment



# Metabolomics Interpretation

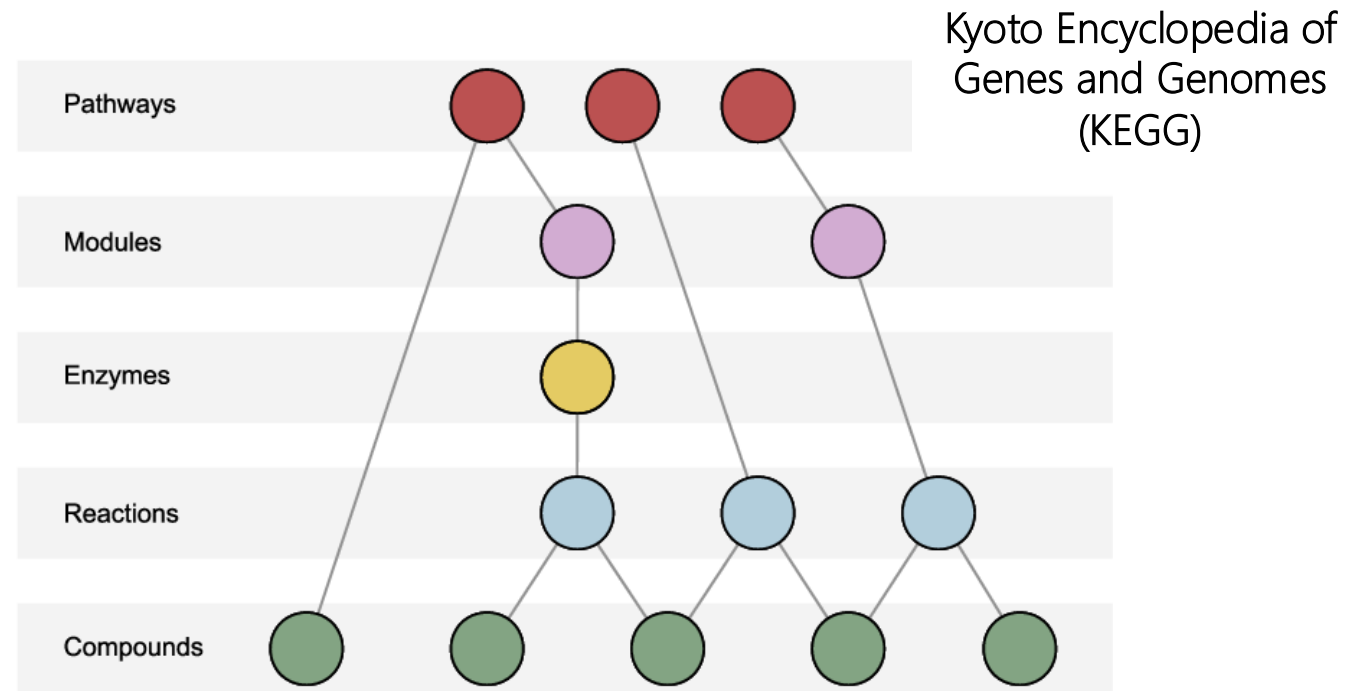
## Over-Representation Analysis (ORA)

Test if the proportion of **differentially expressed molecules** within a particular pathway is statistically higher or lower than would be **expected by chance**.

# Metabolomics Interpretation

## Over-Representation Analysis (ORA)

Test if the proportion of **differentially expressed molecules** within a particular pathway is statistically higher or lower than would be **expected by chance**.



Picart-Armada, Sergio, et al. "FELLA: an R package to enrich metabolomics data." BMC bioinformatics

# Metabolomics Interpretation

## Over-Representation Analysis (ORA)

Background: 100 metabolites measured (that belong to 3 pathways)

Pathway	Total # of metabolites
Pathway A	10
Pathway B	20
Pathway C	30
Not in a pathway	40
Total	100

# Metabolomics Interpretation

## Over-Representation Analysis (ORA)

Background: 100 metabolites measured (that belong to 3 pathways)

10 differentially expressed metabolites

Pathway	Total # of metabolites	Significant metabolites
Pathway A	10	5
Pathway B	20	2
Pathway C	30	1
Not in a pathway	40	2
Total	100	10

# Metabolomics Interpretation

## Over-Representation Analysis (ORA)

Background: 100 metabolites measured (that belong to 3 pathways)

10 differentially expressed metabolites

How many metabolites would we expect by chance?

Pathway	Total # of metabolites	Significant metabolites	Random
Pathway A	10	5	$10 \times (10/100) = 1$
Pathway B	20	2	$10 \times (20/100) = 2$
Pathway C	30	1	$10 \times (30/100) = 3$
Not in a pathway	40	2	
Total	100	10	

# Metabolomics Interpretation

## Over-Representation Analysis (ORA)

Background: 100 metabolites measured (that belong to 3 pathways)

10 differentially expressed metabolites

How many metabolites would we expect by chance?

Pathway	Total # of metabolites	Significant metabolites	Random
Pathway A	10	5	$10 \times (10/100) = 1$
Pathway B	20	2	$10 \times (20/100) = 2$
Pathway C	30	1	$10 \times (30/100) = 3$
Not in a pathway	40	2	
Total	100	10	

# Metabolomics Interpretation

Over-Representation Analysis (ORA)

Contingency table for Pathway A

	In A	Not in A	Total
Significant	5	5	10
Not significant	5	85	90
Total	10	90	100


Hypergeometric distribution:  $P(X \geq 5)$



**Break!**

# Interpretation – Hands on

# Acore enrichment functionality



Search  \* + K

Acore Documentation

API usage examples

- Metabolomics data filtering
- Normalization of samples
- Batch correction of samples
- Imputation data (MS-example)
- Imputation (Metabolomics example)
- Exploratory Analysis
- Metabolomics drift correction
- Differential regulation: T-test for two groups
- Differential regulation: ANOVA across more than two groups and posthoc t-tests
- Differential regulation: ANCOVA for two groups
- Enrichment analysis**



## Enrichment analysis

Is done separately for up- and downregulated genes as it's assumed that biological processes are regulated in one direction.

► Show code cell source

	identifier	group1	group2	pvalue	padj	rejected	log2FC	FC
98	O43432	deceased	alive	0.000	0.001	True	-0.741	0.598
32	Q96JJ3	deceased	alive	0.001	0.018	True	-0.553	0.682
99	O43175	deceased	alive	0.000	0.001	True	1.321	2.498
97	P39059	deceased	alive	0.000	0.000	True	1.600	3.031

Running the enrichment analysis for the up- and down regulated protein groups separately with the default settings of the function, i.e. a log2 fold change cutoff of 1 and at least 2 protein groups detected in the set of proteins defining the functional annotation.

```
ret = acore.enrichment_analysis.run_up_down_regulation_enrichment(  
    regulation_data=diff_reg,  
    annotation=annotations,  
    pval_col="padj",  
    min_detected_in_set=2,  
    lfc_cutoff=1,  
)
```

No significant enrichment found with the given parameters. Returning an empty DataFrame.



Q Search \* + K

Acore Documentation

#### API usage examples

- Metabolomics data filtering
- Normalization of samples
- Batch correction of samples
- Imputation data (MS-example)
- Imputation (Metabolomics example)

Exploratory Analysis

Metabolomics drift correction

Differential regulation: T-test for two groups

[Differential regulation: ANOVA across more than two groups and posthoc t-tests](#)

Differential regulation: ANCOVA for two groups

Enrichment analysis

#### Reanalysis

[Data Analysis PXD040621](#)



Show code cell source



### Enrichment Analysis regulated vs non-regulated

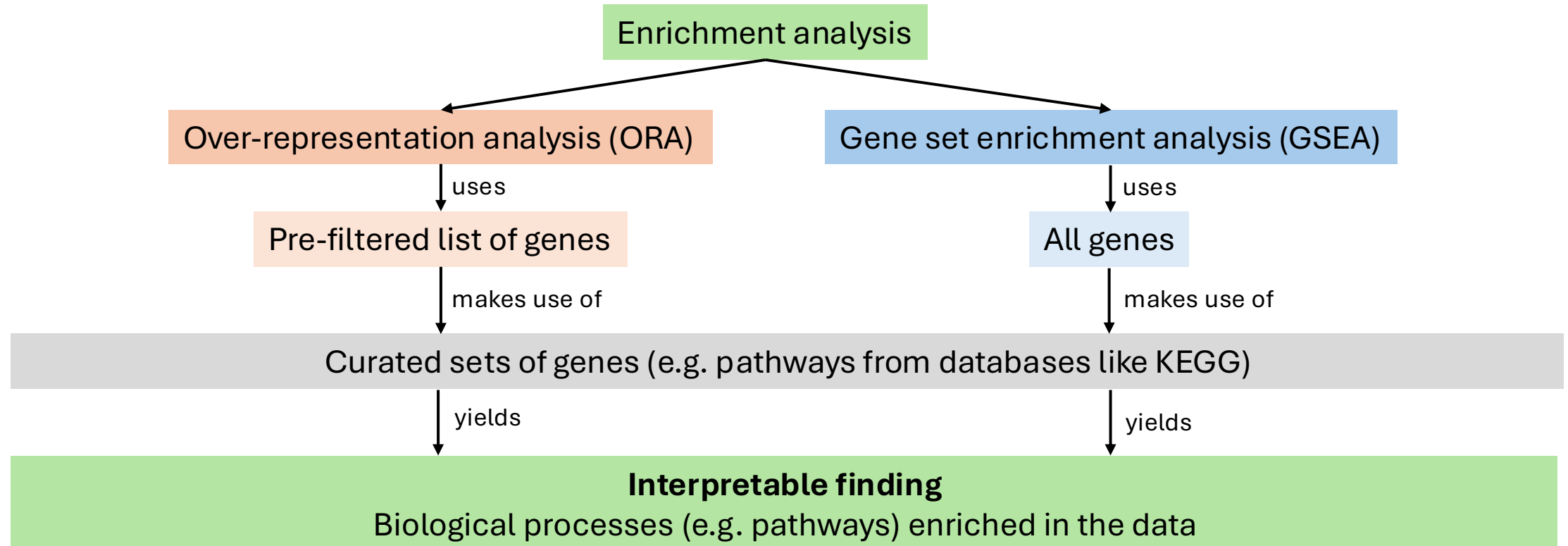


Some slides from  
transcriptomics course

# Functional analysis for biological interpretation

- **Enrichment analysis**

- Knowledge-based approach for biological interpretation of omics data, e.g. RNAseq results
- **Goal:** find biological terms (curated gene sets for biological processes) that are enriched in the data for biological interpretation (no single-gene interpretation)
- Biological term/pathway lists are obtained from KEGG, GO etc or can be manually curated (often .gmt files)
- Different types of enrichment analysis



# Functional analysis for biological interpretation

- **Gene-set enrichment analysis (GSEA)** – all genes ranked by fold-change and p-value

List of **all** genes

	Gene symbol	p-val	Differentially expressed	FC
1	<i>HAO1</i>	0.0001	Yes	-5.6
2	<i>RDH16</i>	0.00045	Yes	-7.6
3	<i>ALDOB</i>	0.00009	Yes	-6.7
4	<i>AKR1C4</i>	0.0051	Yes	10.0
5	<i>ABCB11</i>	0.067	Yes	10.0
6	<i>BAAT</i>	0.0034	Yes	23.5
7	<i>SLC2A2</i>	0.00990	Yes	2
8	<i>FCER1G</i>	0.00673	Yes	6.7
...				
6532	<i>OSM</i>	0.00009	Yes	9
6533	<i>VEGFA</i>	0.067	Yes	8.5
6534	<i>ZAP70</i>	0.0034	Yes	7.8

$$\text{RANKING} = \text{SIGN}(\text{FC}) * -\log_{10}(\text{p-value})$$

Ranked list of genes

	Gene symbol	Ranking
1	<i>JAK9</i>	12.9
2	<i>RDH16</i>	11.4
3	<i>ALDOB</i>	11.3
5	<i>AKR1C4</i>	10.7
6	<i>ABCB11</i>	8.6
7	<i>BAAT</i>	3.6
8	<i>SLC2A2</i>	0.0099
9	<i>FCER1G</i>	0.00673
...		
6532	<i>MyD88</i>	-9
6533	<i>TRIF</i>	-9.37
6534	<i>IRAK1</i>	-15.6

SIGNIFICANT + UPREGULATED

NON-SIGNIFICANT

SIGNIFICANT + DOWNREGULATED



# Functional analysis for biological interpretation

- **Enrichment analysis – details and summary**

	Over-representation analysis (ORA)	Gene set enrichment analysis (GSEA)
Input gene list	Filtered list, e.g. $p_{adj} < 0.05$ and $\log_2FC >  1 $ → Genes not passing the thresholds will be missed	Entire gene list ranked by fold-change (with sign) and/or $p_{adj}$ → No pre-filtering of genes is an advantage
Separation of up- and down-regulated genes	Can be done	No
Core statistical result and tests used	<ul style="list-style-type: none"> <li>• Typically an adjusted p value</li> <li>• Fisher's exact test (one-sided), hypergeometric distribution</li> </ul>	<ul style="list-style-type: none"> <li>• Normalized enrichment score (NES) and adjusted p value</li> <li>• Modified Kolmogorov-Smirnov and permutation test</li> </ul>
Software packages (selection)	mulea, g:profiler, Enrichr, clusterProfiler, GSEApy (via Enrichr)	fgsea, clusterProfiler, GSEApy