

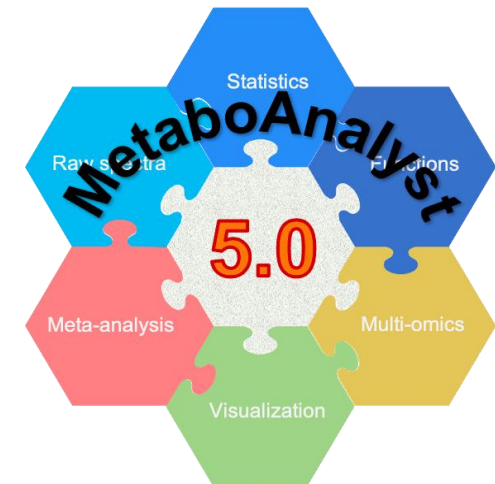
# Spectra processing, functional integration and covariate adjustment of global metabolomics data using MetaboAnalyst 5.0

---

Section I: Metabolomics Raw Spectra data  
processing & functional analysis

TA: Zhiqiang Pang  
([zhiqiang.pang@mail.mcgill.ca](mailto:zhiqiang.pang@mail.mcgill.ca))

18<sup>th</sup> Annual Conference of the Metabolomics Society  
**METABOLOMICS 2022**  
Valencia, Spain | JUNE 19-23  
Pre-Conference Workshops

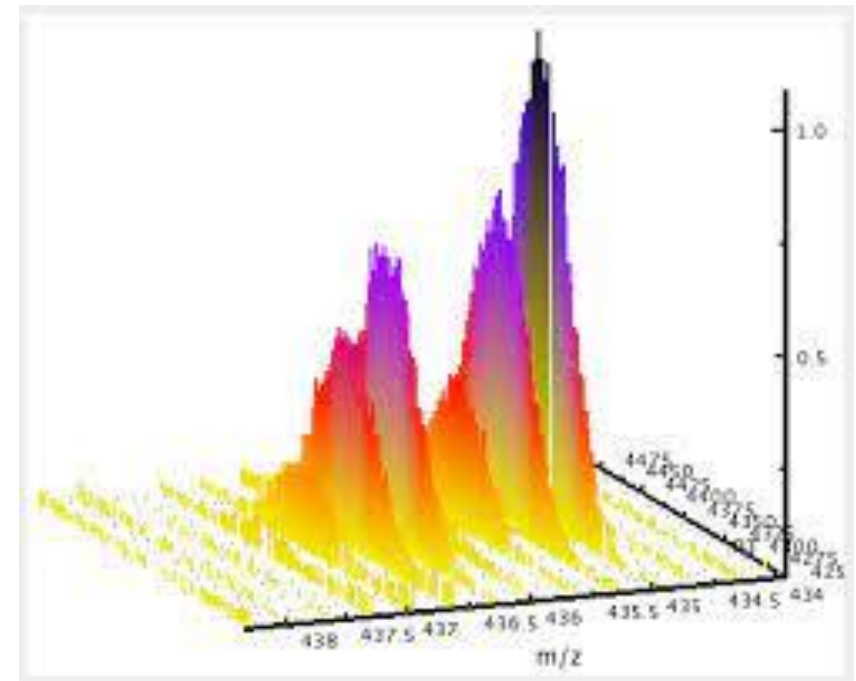


# Learning Questions..

- What the algorithm MetaboAnalyst is using for raw spectra processing?
- How many data files we are supporting for processing online at the maximum?

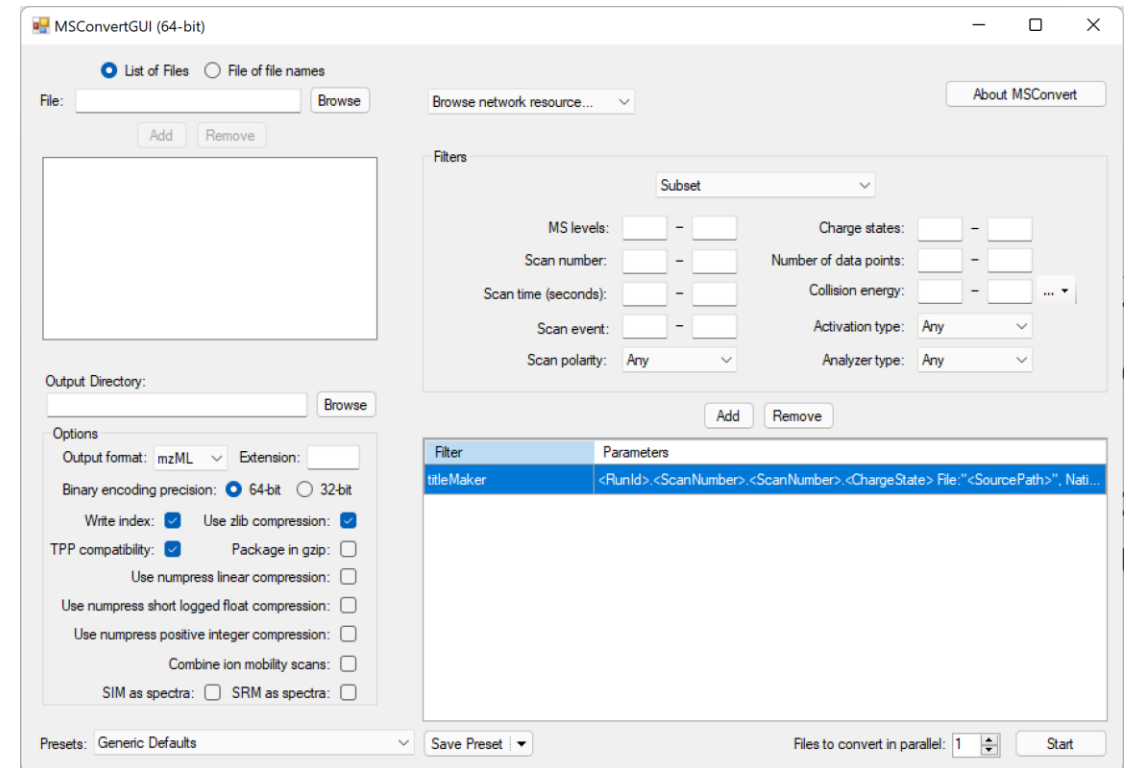
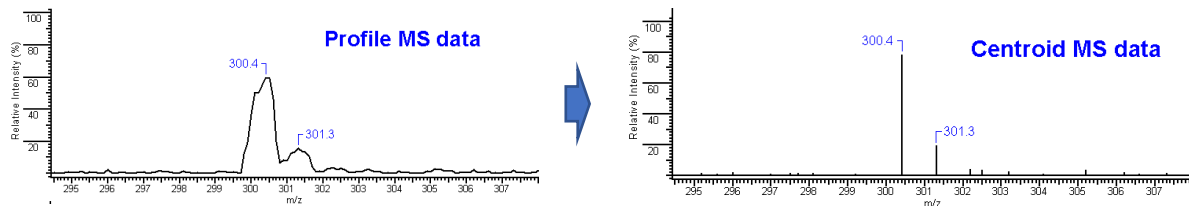
# What is raw spectra (pre-)processing?

- Convert the raw spectra data from MS instrument into metabolic features (MS peaks);
- Usually contains 2~3 dimensions. For DI-MS, the raw spectra includes  $m/z$  and intensity (2D); while for the LC-MS, the raw spectra data includes  $m/z$ , retention time and intensity (3D);
- The most common vendor raw data file formats are .raw/.RAW/.wiff/.d/.D/ etc. They need to be converted into open-source format for further processing with open-source software.

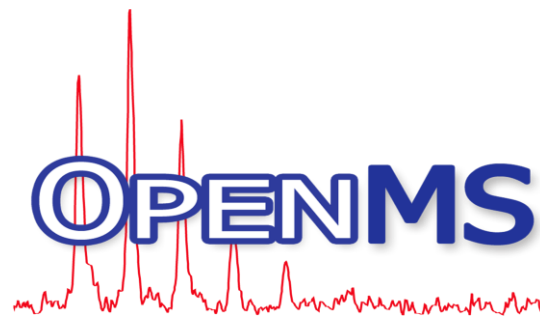
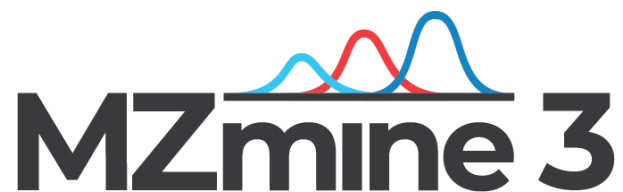
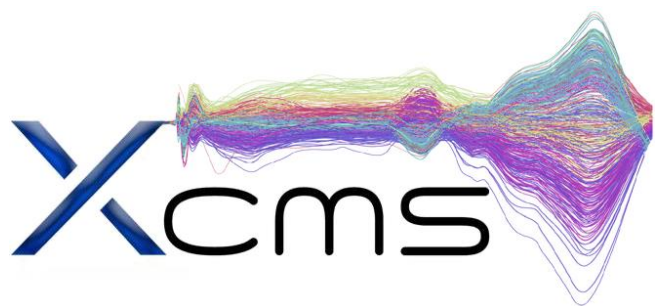


# Profile or Centroid?

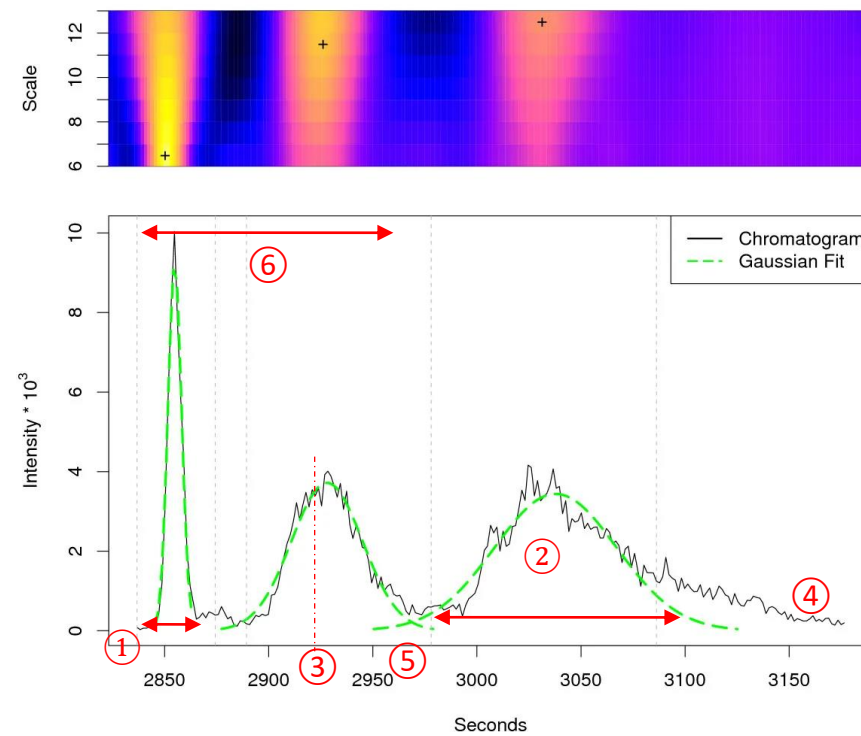
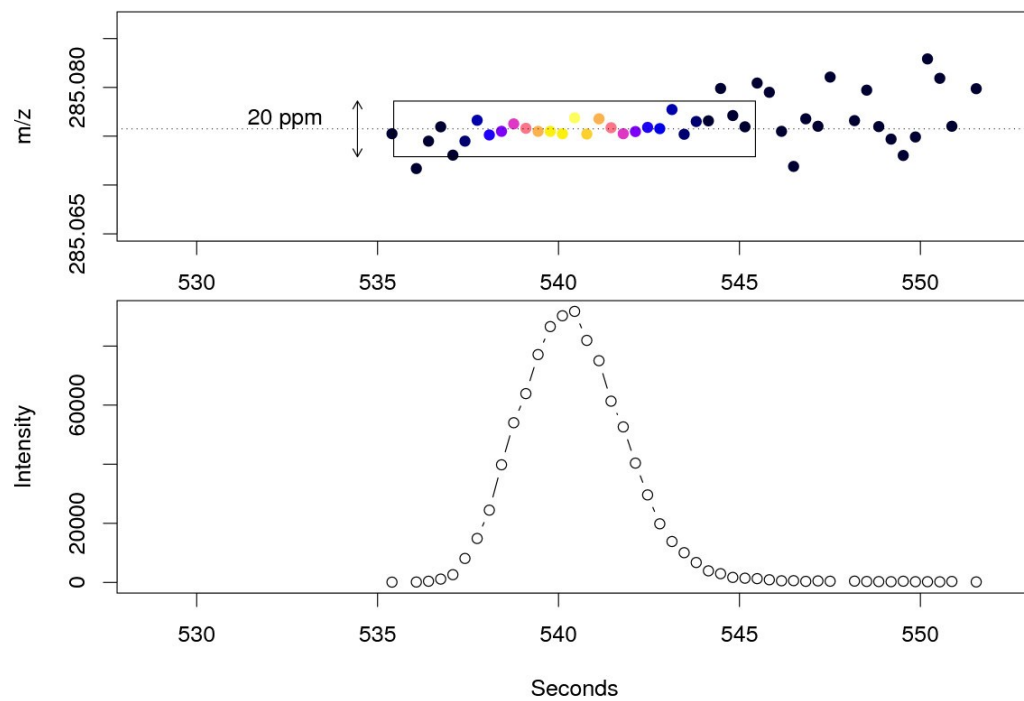
- The vendor raw spectra data is usually in profile format, which is redundant for regular LC-MS based metabolomics analysis;
- We need to convert the MS data into centroid mode to condense the Gaussian Profile peaks into centroids.
- Open-source formats (.mzML/ etc.)..



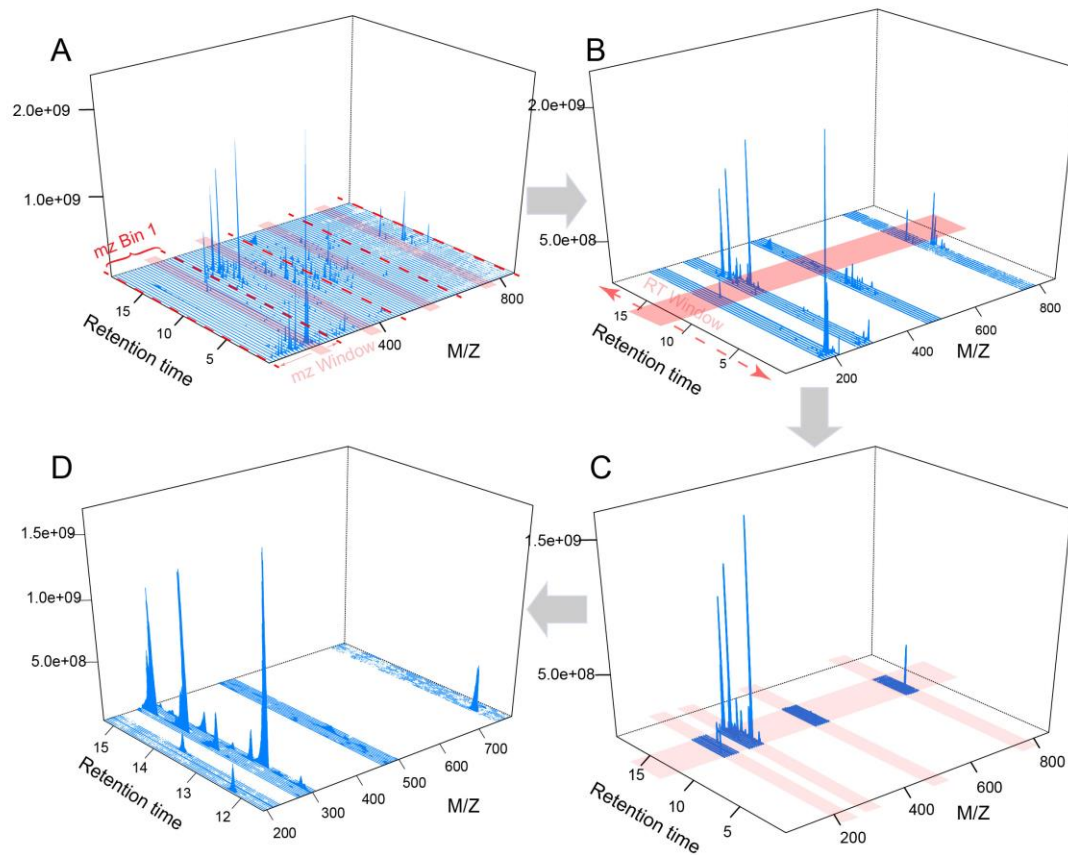
Open-source Software for raw spectra processing..



# centWave



# ROI Extraction



- Data-driven ROI extraction;
- Regions with high abundance of MS signals;
- Both low intensity peaks as well as high intensity peaks will be retained;

# DoE-based Parameter Optimization

DoE -- Central composite design

44 runs

Order	Peakwith_min	Peakwith_max	mzdiff	snthresh	bw
1	-1	-1	-1	-1	-1
2	1	-1	-1	-1	-1
3	-1	1	-1	-1	-1
...	...	...	...	...	...
43	0	0	0	0	1
44	0	0	0	0	0

3 level for every parameters (-1, 0, 1)

Relative reliable peaks ratio  
(identified by their isotopes)

A co-efficient describing the stability of  
a grouped feature

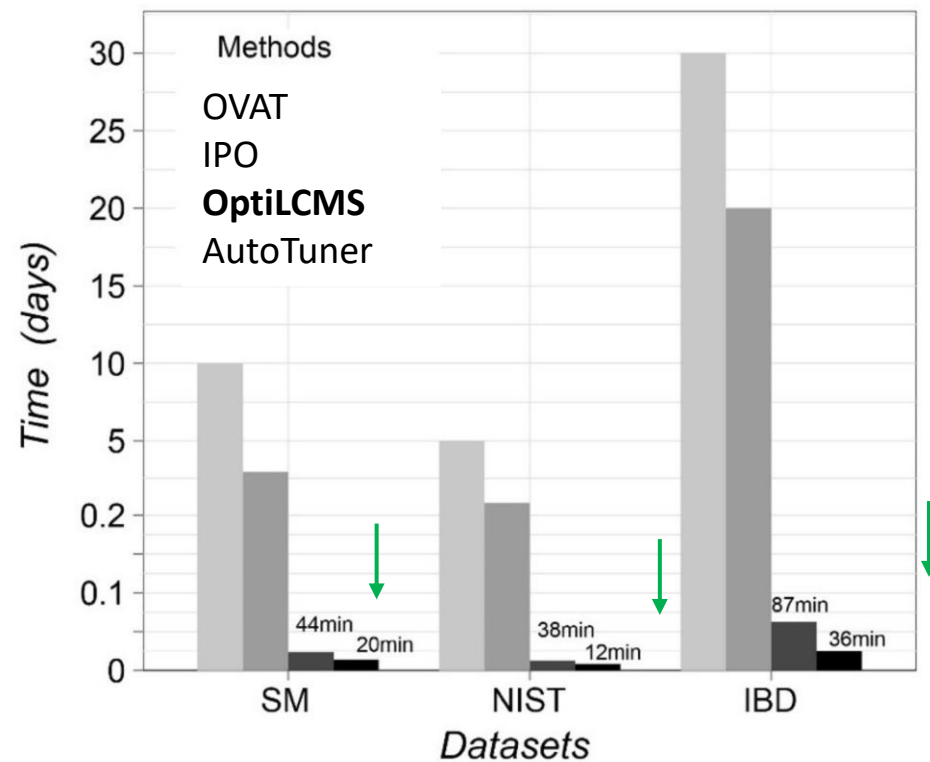
$$QS = \frac{RP^{3/2}}{'all\ peaks' - LIP} * GR^2 * Q_{coE}$$

Gaussian peaks ratio

- The most important parameters are evaluated with 44 DoE runs
- Instead of  $3^8 = 6561$  one-variable-at-a-time runs.



# Performance Evaluation - Speed



+ 3 datasets: Standard Mixture (SM), NIST-SRM 1950 and IBD data from iHMP2.

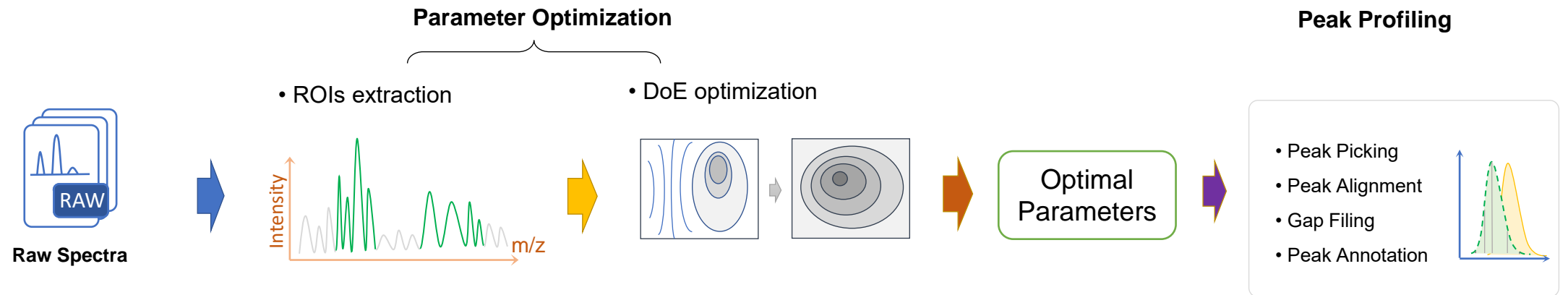
# Bench markings


	Default	Optimized
<b>Total peaks</b>	2,492	2,423
<b>Isotopes / Adducts</b>	667 (26.8%)	1,112 (45.9%)
<b>Formula Assigned</b>	663	762
<b>Potential compounds</b>	1,085	1,692
<b>Variance (PC1 + PC2)</b>	37%	50%
<b>Significant peaks</b>	855	1,091

	Default	Optimized
<b>Total peaks</b>	4,344	5,113 (+ 17.7%)
<b>Isotopes</b>	760	1,274 (+ 67.6%)
<b>Adducts</b>	927	1,132 (+ 22.1%)
<b>Formulas assigned</b>	632	687 (+ 8.7%)
<b>Potential compound matches</b>	1,587	1,803 (+ 13.6%)
<b>Variance explained (PC1 + PC2)</b>	76.5%	81.3% (+ 4.8%)

# Workflow Overview

- Our optimization approach is designed to extract a region abundant with MS signals for a design of experiment (DoE)-based optimization.





[Home](#)

[Data Formats](#)

[Tutorials](#)

[OmicsForum](#)

[APIs](#)

[Update History](#)

[MetaboAnalystR](#)


[Contact](#)


[User Stats](#)


[Publications](#)

[COVID-19 Data](#)

[About](#)







# MetaboAnalyst 5.0

- user-friendly, streamlined metabolomics data analysis

## Module Overview

Input Data Type	Available Modules (click on a module to proceed, or scroll down for more details)					
Raw Spectra (mzML, mzXML or mzData)	LC-MS Spectra Processing					
MS Peaks (peak list or intensity table)			Functional Analysis	Functional Meta-analysis		
Annotated Features (compound list or table)		Enrichment Analysis	Pathway Analysis	Joint-Pathway Analysis	Network Analysis	
Generic Format (.csv or .txt table files)	Statistical Analysis [one factor]	Statistical Analysis [metadata table]	Biomarker Analysis	Statistical Meta-analysis	Power Analysis	Other Utilities

>> Statistical Analysis [one factor]

This module offers various commonly used statistical and machine learning methods including t-tests, ANOVA, PCA, PLS-DA and Orthogonal PLS-DA. It also provides clustering and visualization tools to create dendrograms and heatmaps as well as to classify data based on random forests and SVM.

>> Statistical Analysis [metadata table]

This module aims to detect associations between phenotypes and metabolomics features with considerations of other experimental factors / covariates based on general linear models coupled with PCA and heatmaps for visualization. More options are available for two-factors / time-series data.

>> Biomarker Analysis

This module performs various biomarker analyses based on receiver operating characteristic (ROC) curves for a single or multiple biomarkers using well-established methods. It also allows users to manually specify biomarker models and perform new sample prediction.

Xia Lab @ McGill (last updated 2022-06-17)

# Functional Utilities

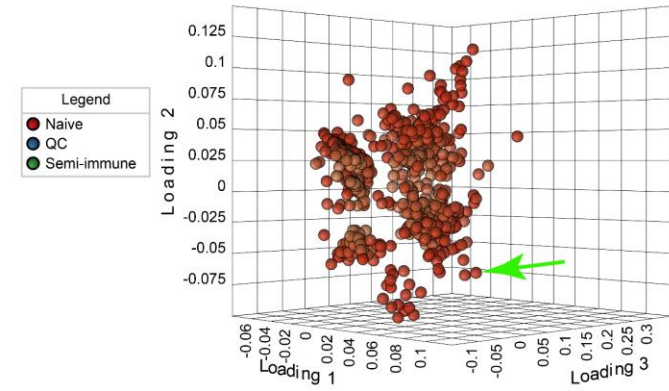
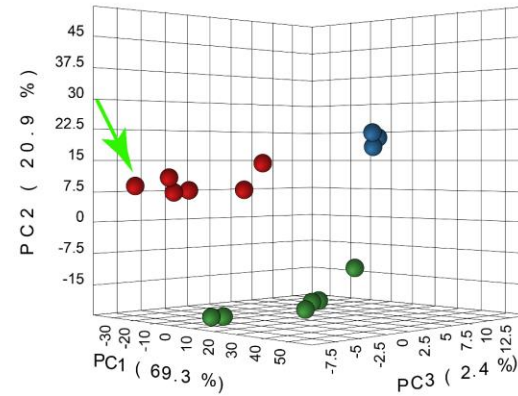
- Raw Data Uploading (.mzML/.mzXML/.mzData/.cdf);
- Centroiding on the fly;
- Parameters Optimization (automatically);
- Peak Profiling (Peak Picking/Alignment/Gap filling);
- Peak Annotation (Adducts + Isotopes);
- Putative Compound Mapping;
- Result visualization...

# Result Demo..

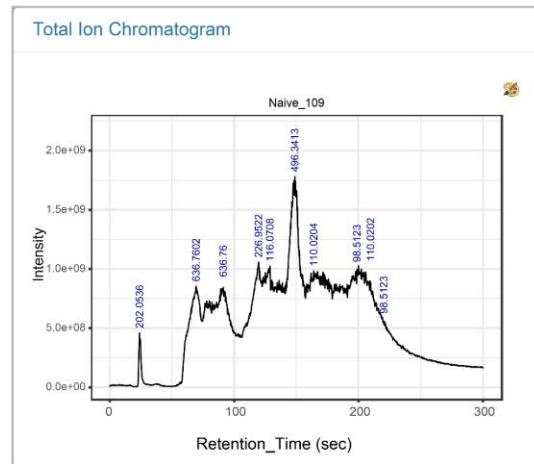
A



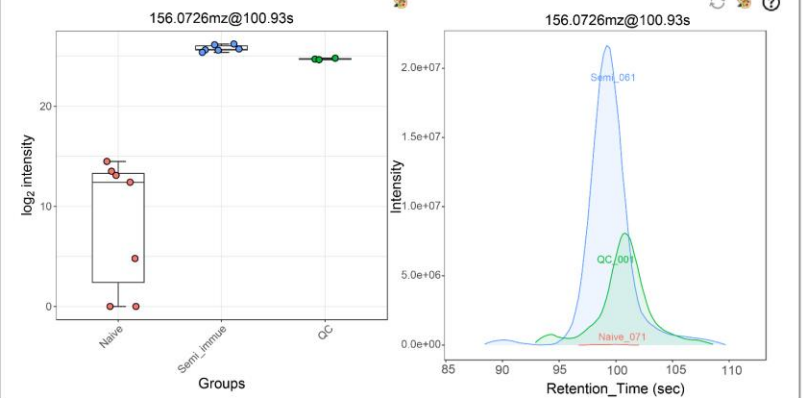
Mouse-drag to rotate, double-clicking a node to view its summary. Use the tables below to [search](#) and [view](#) specific peaks or samples



B



Boxplots and EICs



C

# Functional Analysis

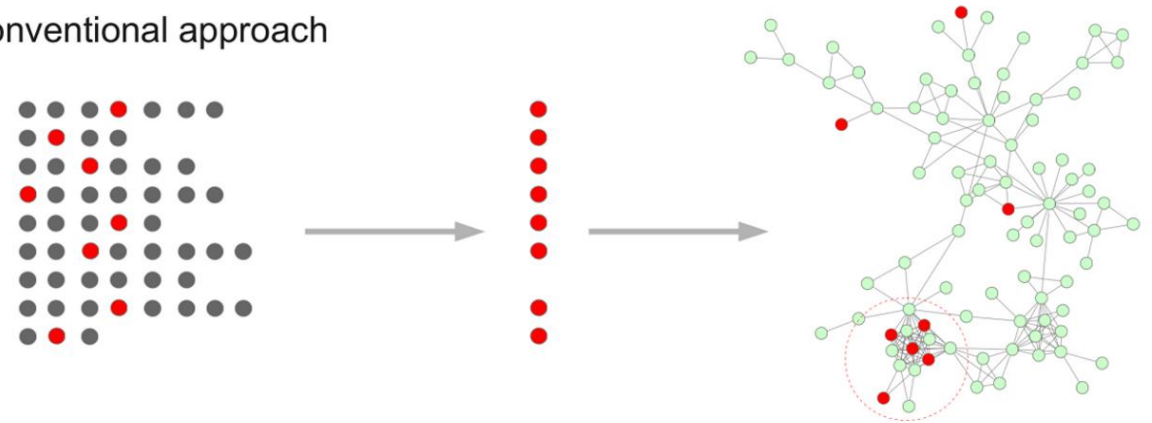
Functions describe the collective behavior of groups of molecules acting together to complete a task (e.g., lysine degradation).

The goal of enrichment analysis is to evaluate whether the members involved in a particular task how more consistent behaviors (e.g., more changes larger than normal) compared with random variations.

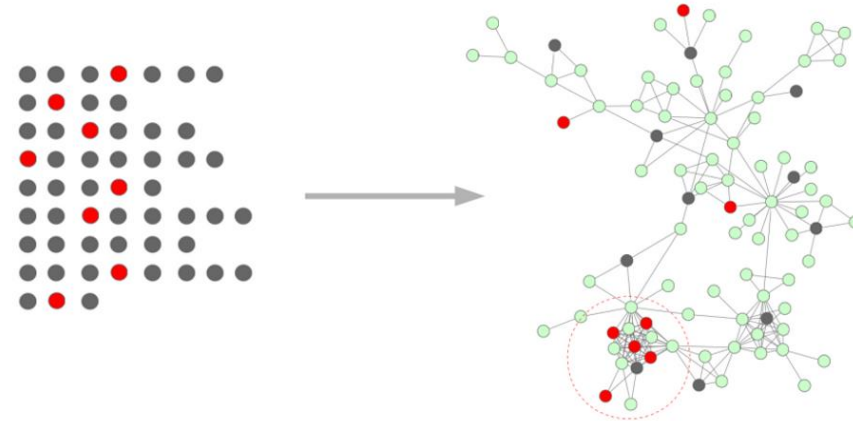
---

***Mummichog*** provides a practical solution of one-step functional analysis, bypassing the bottleneck of upfront metabolite identification.


Conventional approach



**mummichog**



# Implementation in MetaboAnalyst



## MetaboAnalyst 5.0 - user-friendly, streamlined metabolomics data analysis

[Home](#)  
[Data Formats](#)  
[Tutorials](#)  
[OmicsForum](#)  
[APIs](#)  
[Update History](#)  
[MetaboAnalystR](#)  
[Contact](#)  
[User Stats](#)  
[Publications](#)  
[COVID-19 Data](#)  
[About](#)

### Module Overview

Input Data Type	Available Modules (click on a module to proceed, or scroll down for more details)					
Raw Spectra (mzML, mzXML or mzData)	LC-MS Spectra Processing					
MS Peaks (peak list or intensity table)	Functional Analysis		Functional Meta-analysis			
Annotated Features (compound list or table)		Enrichment Analysis	Pathway Analysis	Joint-Pathway Analysis	Network Analysis	
Generic Format (.csv or .txt table files)	Statistical Analysis [one factor]	Statistical Analysis [metadata table]	Biomarker Analysis	Statistical Meta-analysis	Power Analysis	Other Utilities

>> [Statistical Analysis \[one factor\]](#)

This module offers various commonly used statistical and machine learning methods including t-tests, ANOVA, PCA, PLS-DA and Orthogonal PLS-DA. It also provides clustering and visualization tools to create dendrograms and heatmaps as well as to classify data based on random forests and SVM.

>> [Statistical Analysis \[metadata table\]](#)

This module aims to detect associations between phenotypes and metabolomics features with considerations of other experimental factors / covariates based on general linear models coupled with PCA and heatmaps for visualization. More options are available for two-factors / time-series data.

>> [Biomarker Analysis](#)

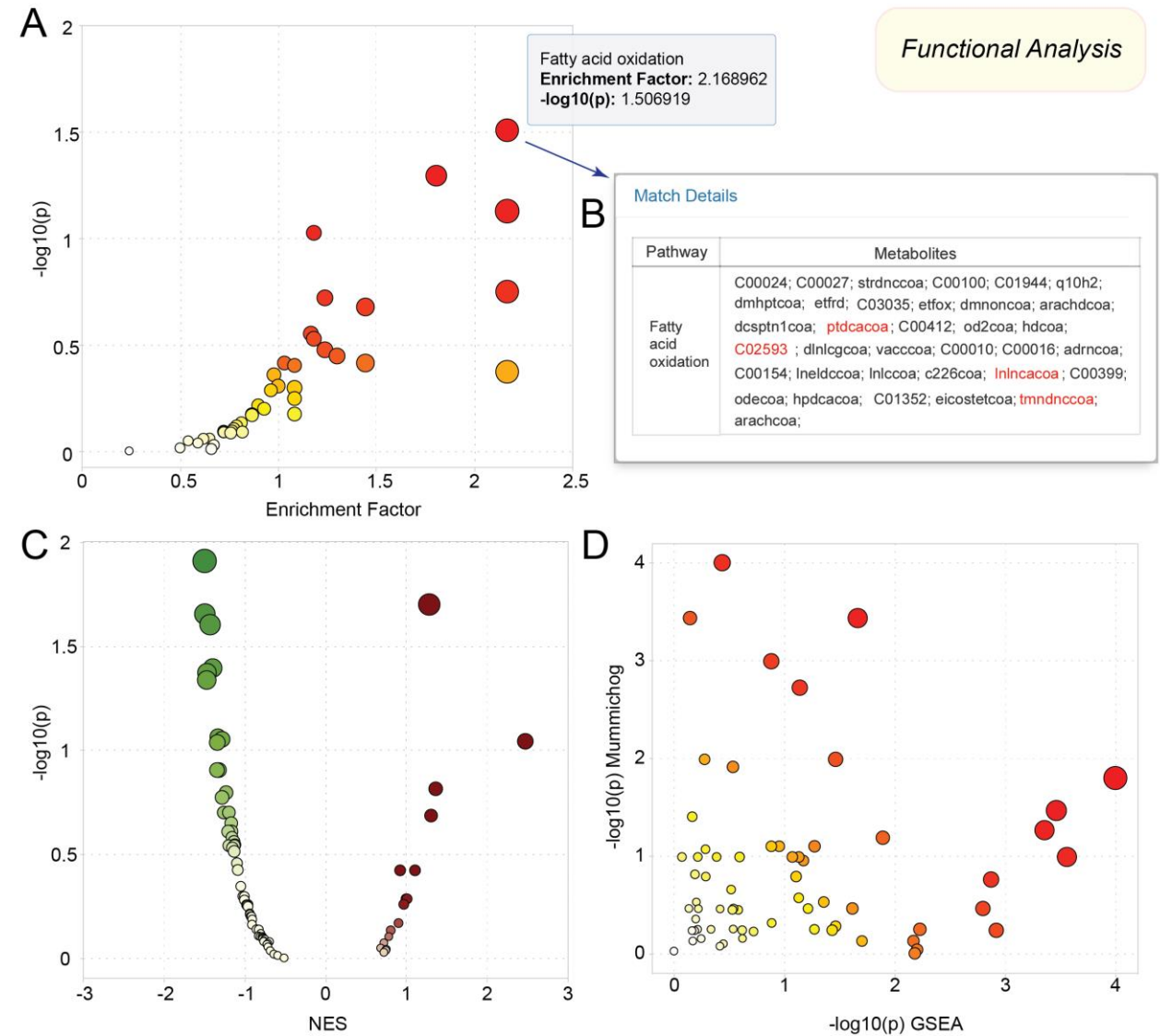
This module performs various biomarker analyses based on receiver operating characteristic (ROC) curves for a single or multiple biomarkers using well-established methods. It also allows users to manually specify biomarker models and perform new sample prediction.

Xia Lab @ McGill (last updated 2022-06-17)

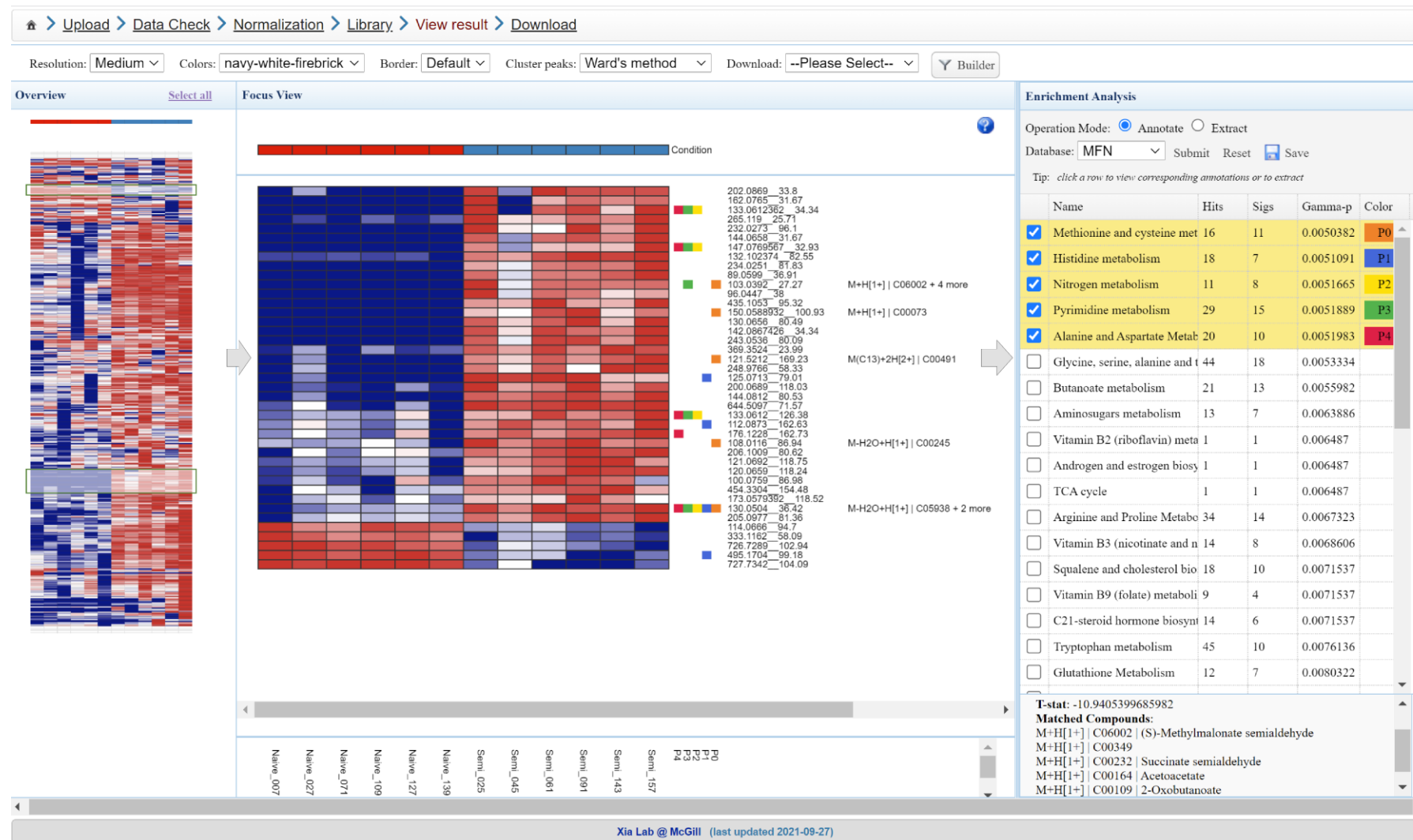


# Functional Analysis Results

- **Mummichog**: implements an over-representation analysis (ORA) method to evaluate pathway-level enrichment based on significant peaks. Users need to specify a pre-defined cutoff.
- **GSEA**: considers the overall ranks of features without using a significance cutoff, and is claimed to be able to detect subtle and consistent changes.
- **Integrated**: combines both Mummichog and GSEA together with Fisher's Method.



# Metabolic Pattern-based functional analysis



# From raw spectra into biological insight

The screenshot displays the MetaboAnalyst web application interface. The browser address bar shows the URL: `dev.metaboanalyst.ca/MetaboAnalyst/Secure/details/DownloadView.xhtml`. The page title is "Download Results & Start New Journey".

**Left Sidebar:** Contains navigation links: Upload, Spectra check, Spectra processing, Job status, Spectra result, Download, and Exit.

**Main Content Area:**

- Download Results & Start New Journey**  
Please download the results (tables and images) from the **Results Download** tab below. The **Download.zip** contains all the files in your home directory. You can also generate a **PDF analysis report** using the button. Finally, you can continue to explore other compatible modules using the **Start New Journey** tab.
- Results Download** (selected) | **Start New Journey**
- General Statistics**
  - ☐ Statistical Analysis [one factor]
  - ☐ Biomarker Analysis
  - ☐ Statistical Analysis [metadata table]
  - ☐ Power Analysis
- Targeted Metabolomics**
  - ☐ Enrichment Analysis
  - ☐ Pathway Analysis
- Global Metabolomics**
  - ☒ Functional Analysis
- GO!** (button, highlighted with a red arrow)

## Tutorials

- Publication: <https://www.nature.com/articles/s41596-022-00710-w>
- Or our manuscript: <https://www.dropbox.com/s/7184c4dheeiiz2p/NP-MetaboAnalyst-2022.pdf?dl=0>
  - ➔ Stage 1: LC–HRMS raw spectra processing
  - ➔ Stage 2: functional analysis of LC–HRMS peaks

## Questions?

- <https://www.omicsforum.ca/>
- If your question is not covered, please create a new topic – we will try to answer them in the coming days.

## Caution:

1. For raw spectra processing, you are strongly encouraged to use 1<sup>st</sup> example rather than the 2<sup>nd</sup> one to avoid waiting in queue for learning purpose;
2. Avoid downloading and uploading any example raw spectra data due the limited bandwidth.
3. Default MetaboAnalyst includes [www.metaboanalyst.ca](http://www.metaboanalyst.ca) and [dev.metaboanalyst.ca](http://dev.metaboanalyst.ca), please use backup node [genap.metaboanalyst.ca](http://genap.metaboanalyst.ca) ONLY if the defaults are not accessible.