M E T A B O L O M I C S

# The X-windows High-Throughput Mass Spectroscopy Pipeline

**Rui C. Martins*,** Castro CC**, Teixeira JA**, Silva-Ferreira AC ***

**\*BioInformatics and Biophysics Lab**
Molecular Biology and Ecology Research Centre
University of Minho
Braga  - Portugal

**\*\*\*Chemiomics and Metabolomics Lab**
Biotechnology Research Center
Portuguese Catholic University
Porto  - Portugal

**\*\*BioEngineering Lab**
Institute for BioEngineering and BioTechnology
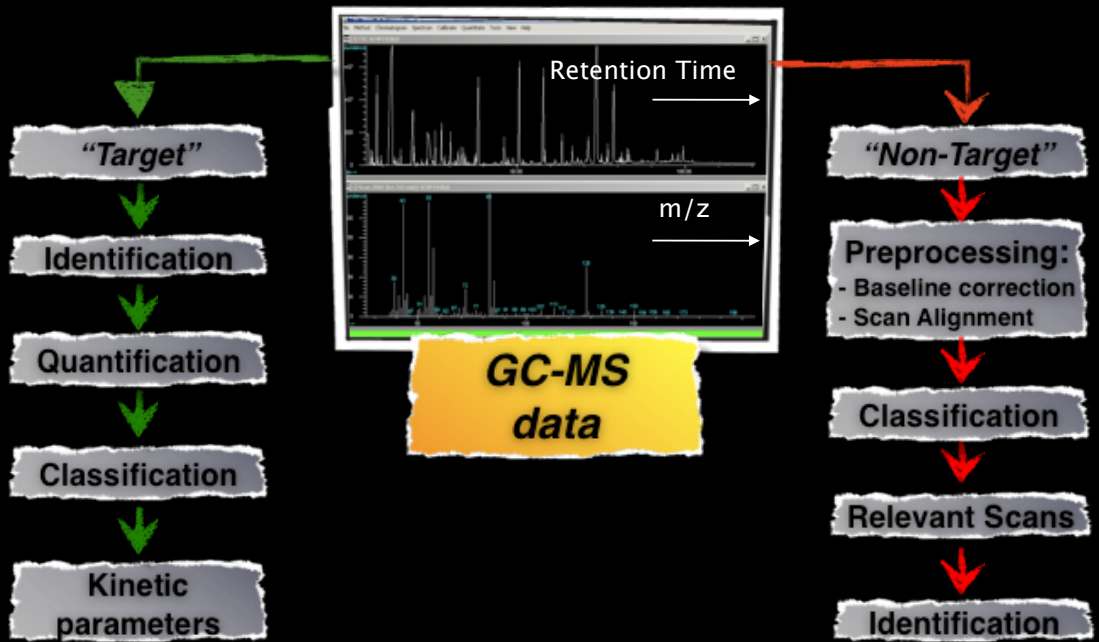University of Minho
Braga  - Portugal

**In This Talk:**

1. Metabolomics: Target vs Non-Target Approaches

2. Mass Spectroscopy Signal Processing

3. The X – Metabolomics Architecture

4. Software Features

5. High-throughput examples

6. Software Demonstration

METABOLOMICS
High-Throughput Mass Spectroscopy

1. Metabolomics: Target vs Non-Target Approaches
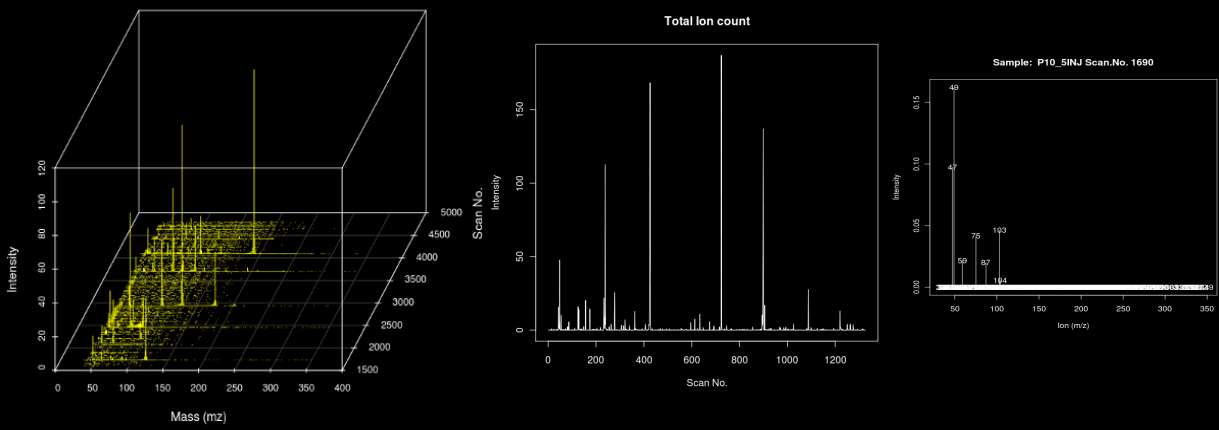
**2. Mass Spectroscopy Signal Processing**

2.1. The MS Chromatographic Signal
2.2. Main Approaches: BioInformatics Vs Chemometrics
2.3. Pre-Processing
2.4. Feature Extraction
2.4. Chromatographic Alignment
2.5. Robust Peak Recognition
2.6. Identification and Composition
2.7. High-troughput MS BioInformatics

METABOLOMICS
High-Throughput Mass Spectroscopy

**2. Mass Spectroscopy Signal Processing**
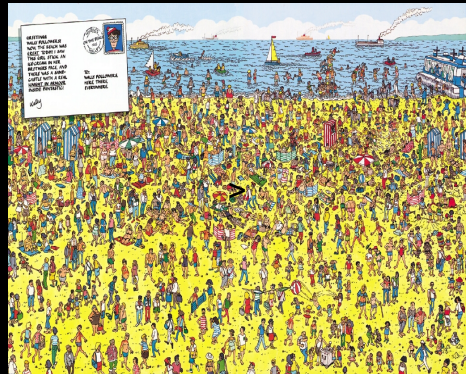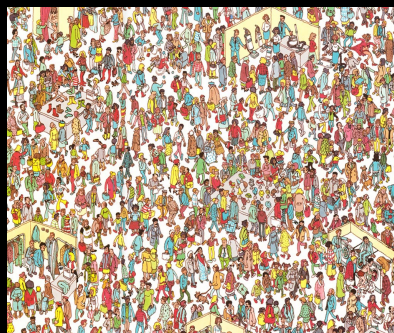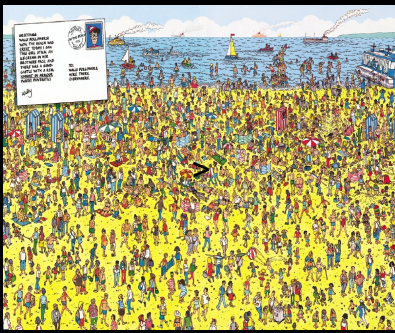
**2.1. The MS Chromatographic Signals**

## 2. Mass Spectroscopy Signal Processing

### 2.1. The MS Chromatographic Signals

**Where Is Wally?**
**or**
**How Can I Know Everybody?**

## 2. Mass Spectroscopy Signal Processing

### 2.1. The MS Chromatographic Signals



**Target Approach**



**Holistic Approach**

**2. Mass Spectroscopy Signal Processing**

**2.1. The MS Chromatographic Signals**

**High-Throughput MS Signal Processing** → **Holistic Approach Complex Systems Systems Biology Systems Chemistry**

## 2. Mass Spectroscopy Signal Processing

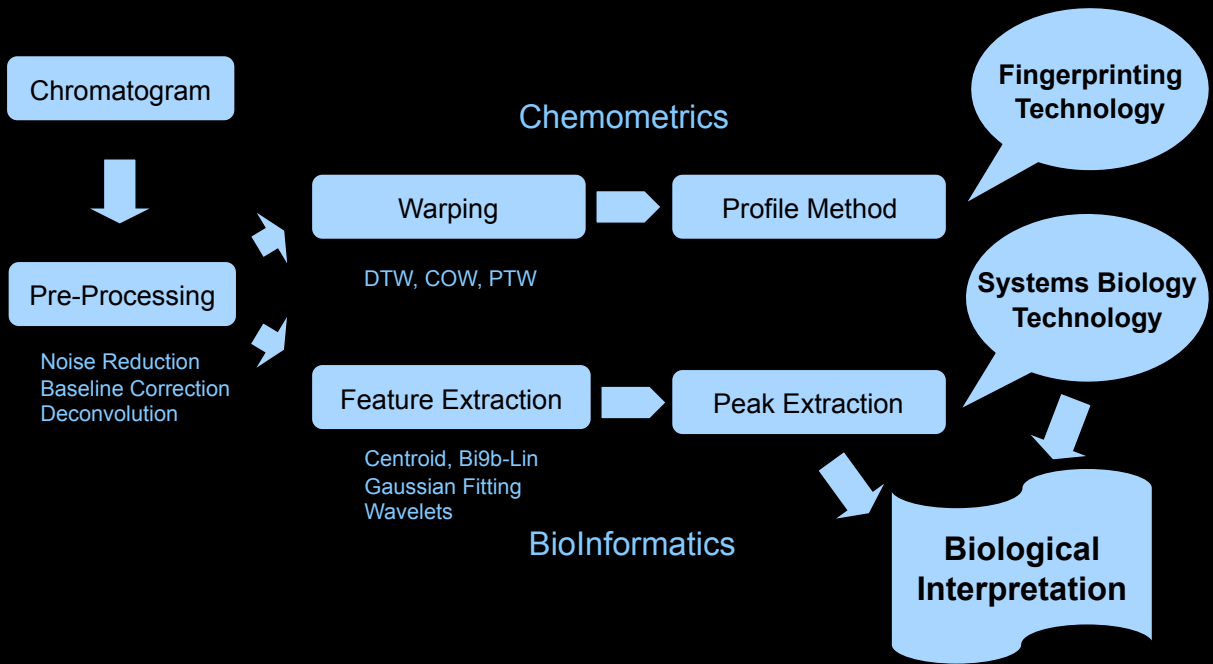### 2.1. The MS Chromatographic Signals

Sample A

Sample B

Sample C

**'In-Silico' Processing:**
Pre-process, Extract, Recognize, Identify, Quantify, Analise
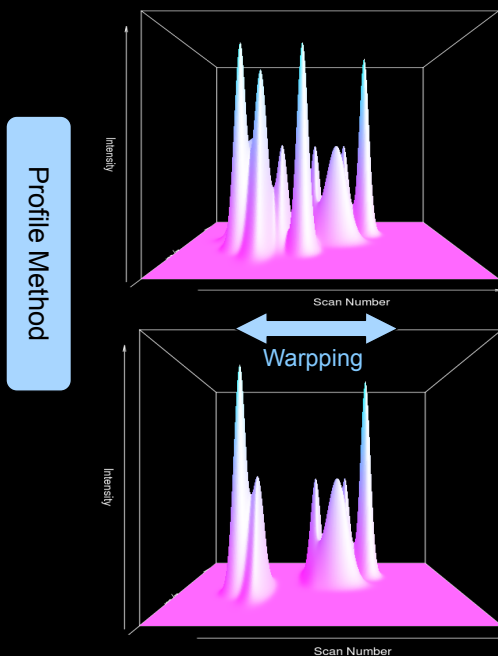
## 2. Mass Spectroscopy Signal Processing

### 2.2. Main Approaches: BioInformatics Vs Chemometrics

Chromatogram

Pre-Processing

Noise Reduction
Baseline Correction
Deconvolution

Chemometrics

Warping

DTW, COW, PTW

Profile Method

Feature Extraction

Centroid, Bi9b-Lin
Gaussian Fitting
Wavelets

BioInformatics

Peak Extraction

Fingerprinting
Technology

Systems Biology
Technology

**Biological
Interpretation**

## 2. Mass Spectroscopy Signal Processing

2.2. Main Approaches: Chemometrics

Profile Method



Warpping

Warp the chromatogram repecting an objective function:

1.Dynamic Time Warpping:  $j = \text{argmin} \, (yi - yr)^2$

2. Correlation Time Warping: $j = \text{argmax} \, \text{corr}(Yi/Yr)$

e.g. Which Reference Chromatogram to use with complex biological samples?

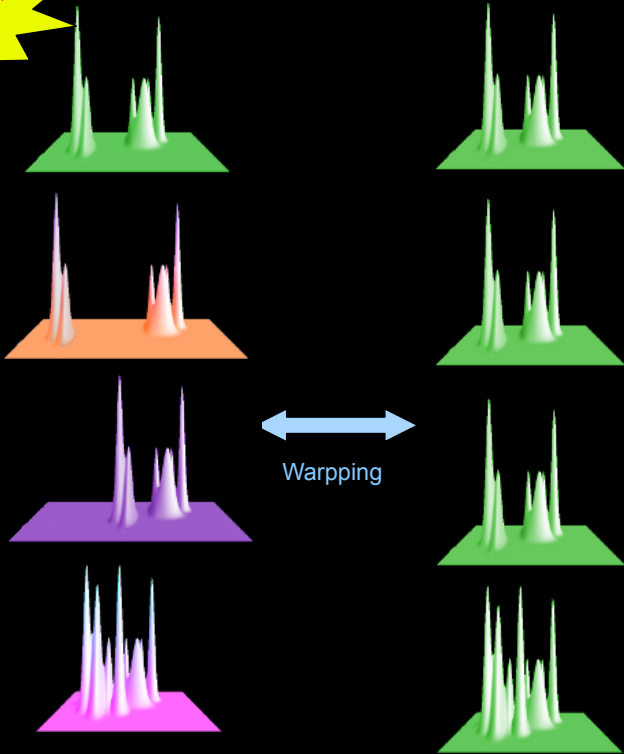**Only robust with low-rank Chromatography – Low complexity of samples → Not for Metabolomics!**

**Not For Complex Chromatograms!**

**2. Mass Spectroscopy Signal Processing**

# 2. Mass Spectroscopy Signal Processing

METABOLOMICS
*High-Throughput Mass Spectroscopy*

Fingerprinting

Profile Method

Sample

Tensor Representation

m/z

Scan No

Unfold Tensor

Fragment No

Samples

Intensity

Multi-Dimensional Linear Algebra

**Multi-Way Chromatogram Decomposition**

Parafac
Tucker 3D
N-way PLS

**Bi-Linear Multivariate Analysis**

SVD
PLS-N
ANN

**METABOLOMICS**
*High-Throughput Mass Spectroscopy*

**2. Mass Spectroscopy Signal Processing**

Fingerprinting

Appropriate for Process Analytical Technology

**Profile Methods**

Not Appropriate For Complex Biological Systems

**Slow Convergence**

**Fingeprinting Only**

**Low Rank Only**

**MS Artifacts**

**Chrom. Artifacts**

**Difficult to provide High-Throughput Metabolomic Information for Systems Biology/Chemistry and Complex Systems Approaches**

**Feature Extraction**

**METABOLOMICS**
*High-Throughput Mass Spectroscopy*

2. **Mass Spectroscopy Signal Processing**

2.4. Feature Extraction

Centroid

Centroid

Gaussian Curve
Peak Extraction

Baseline
Interception

Wavelet
Peak Extraction

Peak Bining

Bin-Lin/Centroid
Methods

## 2. Mass Spectroscopy Signal Processing

### 2.4. Feature Extraction

Wavelets:

- Representation of a signal by a new orthonormal space basis given by non-stationary oscillating waveforms;

- Discontinuities and sharp peaks;

- The Mexican Hat Wavelet:

$$\psi(t) = \frac{1}{\sqrt{2\pi}\sigma^3}\left(1 - \frac{t^2}{\sigma^2}\right)e^{\frac{-t^2}{2\sigma^2}}$$

Wavelet
Peak Extraction

**Frequency**

**Time**

**Multi-Scale Chromatogram Decomposition**

## 2. Mass Spectroscopy Signal Processing

### 2.4. Feature Extraction

Synchronization

Centroid Retention Time
Correction Among Samples

**Bin-Lin/Centroid Methods**

Interpolation functions
To correct deviation in
'non-common' centroids

## 2. Mass Spectroscopy Signal Processing

2.4. Feature Extraction

**Target Tech**

**Bin-Lin/Centroid Methods**
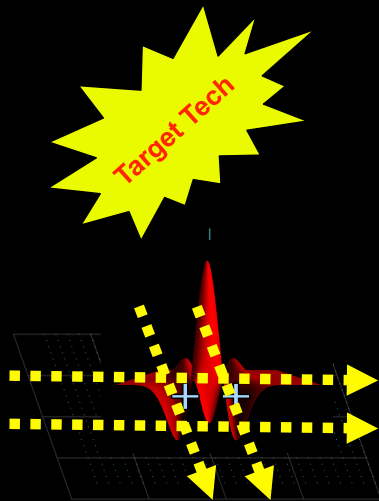
MZMine
MetAlign
XCMS
OpenMS

**Fast Algorithm**

**Simple Peaks**

**Simple Chromatograms**

**Appropriate for**:
Proteomics
LC-MS
Quadrupole

**Fails in:**
GC-MS
Ion-Trap
Complex Chromatograms

**Targeted MS Spectroscopy**

**Difficult to provide High-Throughput Metabolomic Information for Systems Biology/Chemistry and Complex Systems Approaches**

Complex Mathematics

**Feature Extraction**

Higher Efficiency

**2. Mass Spectroscopy Signal Processing**

2.4. Feature Extraction

Target Tech

Failure
In Compound Grouping
(e.g. Loss of Compounds)

Bin-Lin/Centroid
Methods

Failure in Multiscale
Extraction
(e.g. loss of small
peaks)

Failure
In Chromatrogram
Reconstruction
(e.g. Colapse of adjacent
Scans, loss of fragments)

# 2. Mass Spectroscopy Signal Processing

## 2.4. Feature Extraction

How to provide
High-Throughput
Metabolomic
Information?

Signal Processing
Feature Extraction

**Standards for 'In-silico'
High-Throughput Chromatography
In Metabolomics**



High Quality Signals
Peaks High Definition
High Peak Concentration
Filters
Deconvolution
Recognize,
Identify, Quantify,
Analise, Quality Control
Biological Interpretation

**2. Mass Spectroscopy Signal Processing**

2.3. Pre-Processing

METABOLOMICS
*High-Throughput Mass Spectroscopy*

Noise
Supression

Baseline
Correction

Original
Data

M E T A B O L O M I C S

High Quality
Chromatogram

**2. Mass Spectroscopy Signal Processing**

2.3. Pre-Processing

METABOLOMICS
High-Throughput Mass Spectroscopy

**2. Mass Spectroscopy Signal Processing**

2.3. Pre-Processing

Feature Randomization
During on-Trap
Saturation

Deconvolution
Is not possible

Regular Peak
(e.g. phenylethanol)

Saturated Peak
(e.g. phenylethanol)

**Saturation**

**Saturation
Filter**

## 2. Mass Spectroscopy Signal Processing

### 2.3. Pre-Processing

Fuzzy Filtering using Sinkhorn Factorization!!!

Peak = $D_1AD_2$ (perform until convergence)

- Inside a sample and Between samples!!!

Pass

Reject

**Feature Self-Consistency Test**

**Only Robust Peaks Pass**

## 2. Mass Spectroscopy Signal Processing

### 2.3. Pre-Processing

Correlation Between
fragments to obtain
The clusters of
Non-biological
components

Check each cluster
According to
Tikunov et al (2005)

Non-biological
components

**Clean Non-Biological Components**

**Robust
Databases**

**3. The X-Metabolomics Architecture**

3.1. Operating Systems
3.2. Software Compatibility
3.3. File Formats
3.4. Supported Equipments
3.5. Software Architecture in a Nutshell

*In*: *Castro, CC, Teixeira, JA, Silva-Ferreira, AC, Martins, RC. 2009. X-Metabolomics: A high-throughput GC-MS metabolomics pipeline for Saccharomyces cerevisiae. BMC BioInformatics, Submitted.*

METABOLOMICS
High-Throughput Mass Spectroscopy

**METABOLOMICS**

*High-Throughput Mass Spectroscopy*

**3. The X-Metabolomics Architecture**

3.1. Operating Systems - UNIX Like Platforms

Currently Available  for:





**Ubuntu Linux:**
8.04 (Discont.)
8.10 (Stable)
9.04 (Stable)
9.10 (Testing)

**Mac OS X**
Snow Leopard

**METABOLOMICS**
*High-Throughput Mass Spectroscopy*

**3. The X-Metabolomics Architecture**

**3.2. Software Compatibility**

Xorg

R-project

BioInformatics Platforms

OpenPAT Plug-In

MZmine

MetAlign

Markup Languages and DataBases

Mass Spectroscopy Processing Software

METABOLOMICS

High-Throughput Mass Spectroscopy

**3. The X-Metabolomics Architecture**

3.3. File Formats

ASCII Text Files

MZ XML (Proteomics)

NetCDF (Preferencial!!!)
(Use Mass Transit for conversion)

netCDF

NetCDF (network Common Data Form)
Multi-Dimensional Array-oriented Scientific Data
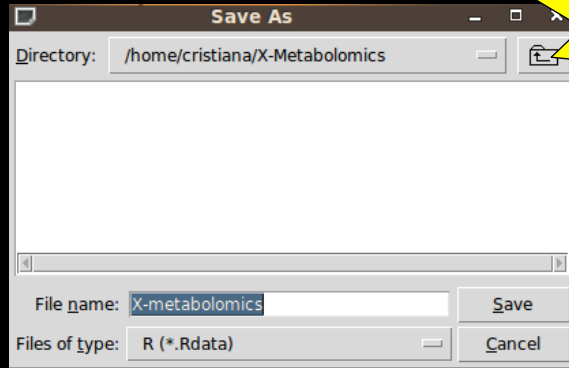Useful for Complex MS/MS datasets

**METABOLOMICS**
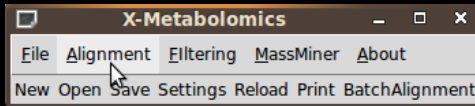*High-Throughput Mass Spectroscopy*

## 4. Software Features

# 4. Software Features

## 4.1. Configuration

**X-Metabolomics**

File  Alignment  FIltering  MassMiner  About

New  Open  Save  Settings  Reload  Print  BatchAlignment

```
>>  ShowXMetabolomicsSettings()
[1] "You are trying to use ShowXMetabolomicsSettings()"
              BinLin CentWave Group
fwhn            1000      NaN   NaN
mzdiff             1      NaN   NaN
ppm              nan      4.0   NaN
snthresh         nan      0.1   NaN
peakwidthmin     nan     15.0   NaN
peakwidthmax     nan     60.0   NaN
scanrangemin     nan    100.0   NaN
scanrangemax     nan   4000.0   NaN
bw               nan      NaN  1.00
minfrac          nan      NaN  0.10
minsamp          nan      NaN  2.00
mzwid            nan      NaN  0.01
max              nan      NaN 50.00
sleep            nan      NaN  0.00
Flag3D         FALSE      NaN   NaN
Filter          0.01      2.0  4.00
```

**Settings: Unix Style**

**Control Feature Extraction and Filters**

# 4. Software Features

## 4.2. File Import



METABOLOMICS
High-Throughput Mass Spectroscopy

# 4. Software Features

## 4.2. File Export
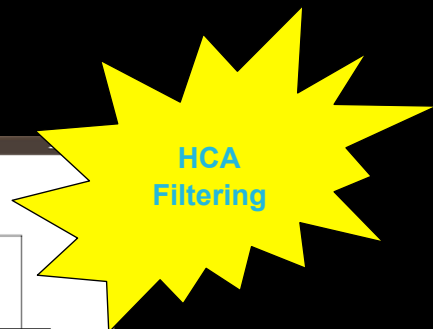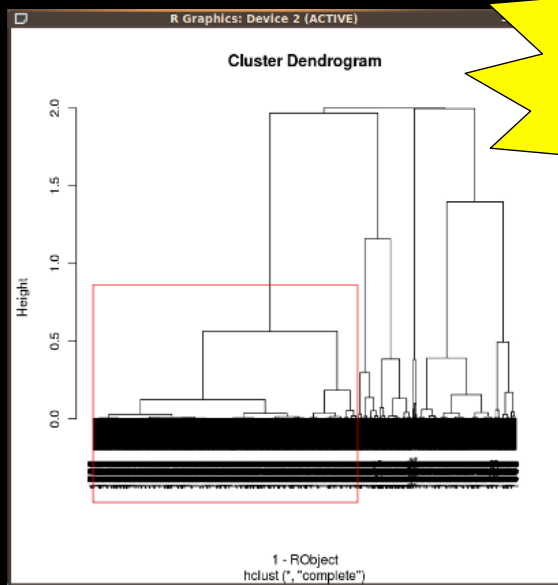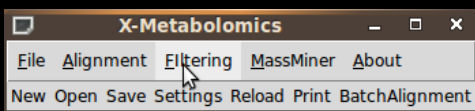
## 4. Software Features

### 4.3. Pre-processing
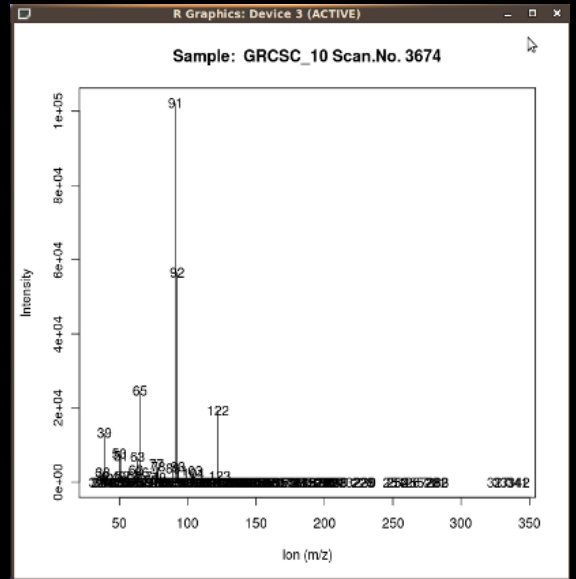
# 4. Software Features

## 4.5. Fingerprinting Diagnostics

# 4. Software Features

## 4.6. MS Quality Control Charts

# 4. Software Features

## 4.7. Unsupervised Metabolomics

4. Software Features

4.9. Co-Expression Pathway Analysis

# 5. Benchmarks

## Minimum Media *Saccharomyces Cerevisiae* Growth



**MinimalMediumFermentation**

Legend:
- X–Met–Centroid–Fast
- X–Met–Centroid–Semiauto.
- X–Met–Wavelet–Fast
- X–Met–Wavelet–Semiauto
- MetAlign
- MZMine

In: Castro, CC, Silva-Ferreira, AC,,Teixeira, JA, Martins, RC. 2009. X-Metabolomics: A High-throughput GC-MS metabolomics pipeline for Saccharomyces cerevisiae . BMC BioInformatics, Submitted.

METABOLOMICS
High-Throughput Mass Spectroscopy

## 5. Benchmarks

Wine Fermentation



*In*: *Castro, CC, Silva-Ferreira, AC,,Teixeira, JA, Martins, RC. 2009. X-Metabolomics: A High-throughput GC-MS metabolomics pipeline for Saccharomyces cerevisiae . BMC BioInformatics, Submitted.*

## 5. Benchmarks

Madeira Wine GC-MS Data
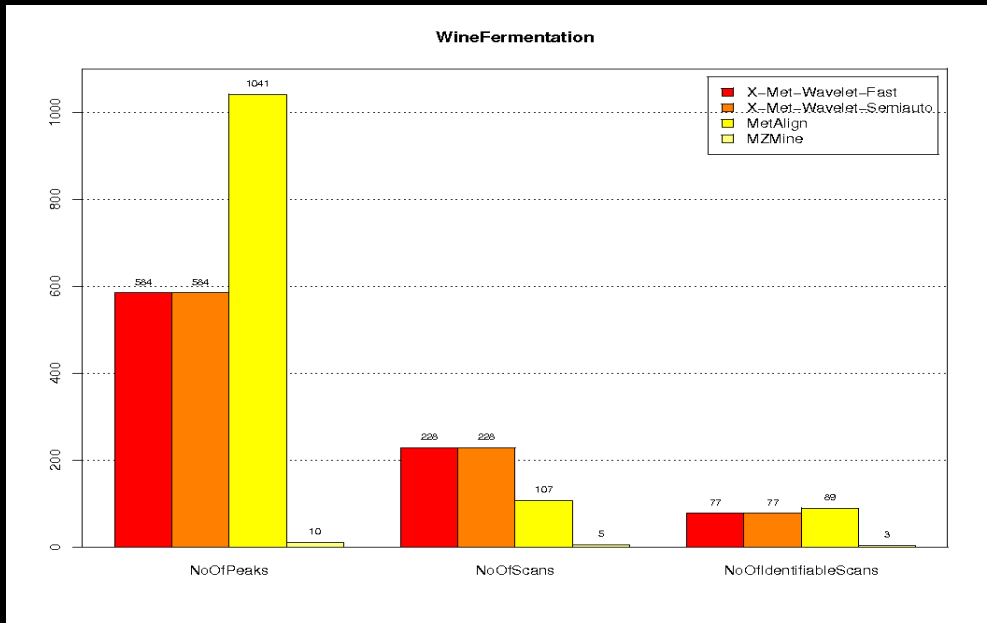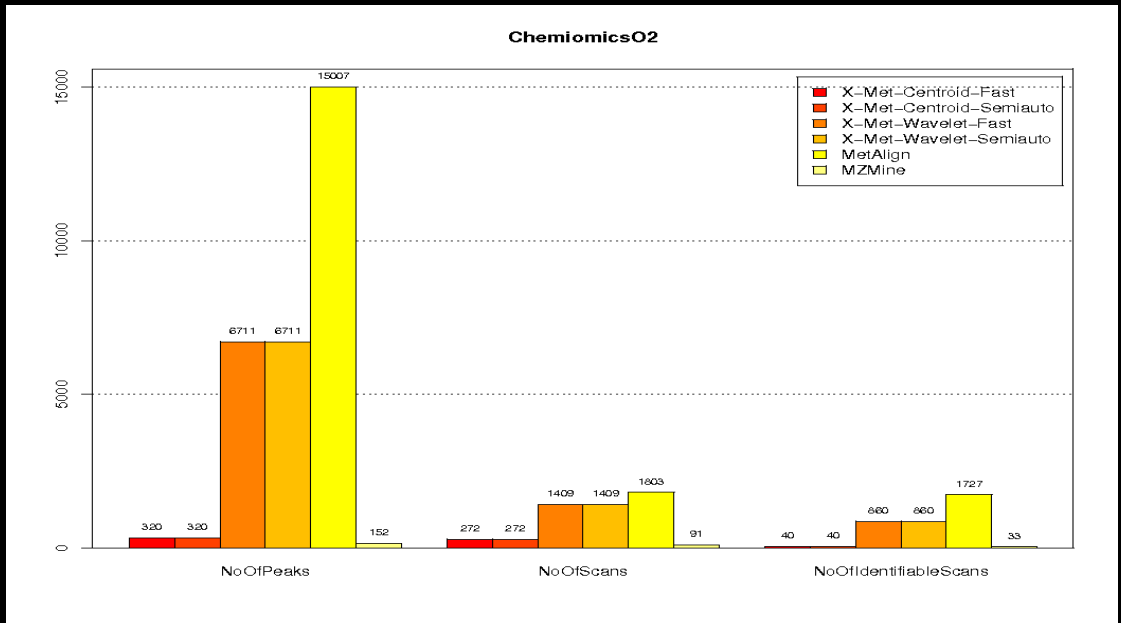


*In: Castro, CC, Silva-Ferreira, AC,,Teixeira, JA, Martins, RC. 2009. X-Metabolomics: A High-throughput GC-MS metabolomics pipeline for Saccharomyces cerevisiae . BMC BioInformatics, Submitted.*

**5. Benchmarks**

**Conclusions**:

**X-Metabolomics** and **MetAlign** were developed for Metabolomics and Metabonomics research – **Out-perform all other software**, which were mainly **developed for Proteomics**

Performance of **X-Metabolomics** is **in most cases similar** to **MetAlign with well defined peaks**

**X-Metabolomics out-performs MetAlign** when chromatograms exhibit **ion-trap artifacts** (due to X-metabolomics filters)

**Further developments in «In-Silico» Chromatography are only possible by the development of new feature extraction methods!**

**METABOLOMICS**
*High-Throughput Mass Spectroscopy*

**6. Software Demonstration**

**Minimum Media Fermentation Time Course Analysis**

*In*: Silva-Ferreira, AC, Gunning, C., Castro, CC, Teixeira, JA,, Martins, RC. 2009. A non-target approach for time-course *Saccharomyces cerevisiae* oxidative response by GC-MS and Cyclic Voltammetry. Metabolomics, Submitted.