

mineXpert2 is part of the msXpertSuite software package  
Modelling, simulating and analyzing ionized flying species

# MINEXP2 User Manual

VISUALIZING, ANALYZING AND MINING OF  $MS^N$  MASS SPECTROMETRIC DATA

---

MINEXP2 7.2.0

# MINEXP2 USER MANUAL: VISUALIZING, ANALYZING AND MINING OF MS<sup>N</sup> MASS SPECTROMETRIC DATA

by Filippo Rusconi

October 15, 2020 , 7.2.0

Copyright 2009-2020 Filippo Rusconi

msXpertSuite - mass spectrometry software suite

[HTTP://WWW.MSXPERTSUITE.ORG/](http://www.msxpertsuite.org/) 


This book is part of the msXpertSuite project.


The msXpertSuite project is the successor of the massXpert project. This project now includes various independent modules:

- massXpert, a program to model polymer chemistries, perform a wide array of (bio-)chemical reactions and simulate the corresponding mass spectrometric data;
- mineXpert2, the successor of mineXpert, a program to visualize and mine MS<sup>n</sup> mass spectral data (mass spectrum, drift spectrum, XIC chromatograms) starting from the TIC chromatogram.

This program is free software: you can redistribute it and/or modify it under the terms of the GNU General Public License as published by the Free Software Foundation, either version 3 of the License, or (at your option) any later version.

This program is distributed in the hope that it will be useful, but WITHOUT ANY WARRANTY; without even the implied warranty of MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE. See the GNU General Public License for more details.

You should have received a copy of the GNU General Public License along with this program. If not, see [HTTP://WWW.GNU.ORG](http://www.gnu.org) ([HTTP://WWW.GNU.ORG/LICENSES/](http://www.gnu.org/licenses/)) .

The flying frog picture is courtesy [HTTP://WWW.PAPUAWEB.ORG](http://www.papuaweb.org) . The specific license as of 20190104 is: *Please acknowledge the use of Papuaweb resources in your publications. To do this include the complete item URL (for example "http://www.papuaweb.org/gb/ref/hinton-1974/63.html") or a general reference to "http://www.papuaweb.org" in your citation/bibliography. This will improve recognition of these resources by Google Scholar and similar search engines.*

## Revision History

Revision 7.0.1	october 2020	fr
Update of the user manual to document the new combine+ and combine- combinations features.		
Revision 7.0.0	september 2020	fr
Update of the user manual to document new features, in particular the MS <sup>n</sup> mass spectrometric data mining and new mass data plotting features.		
Revision 6.0.0	january 2020	fr
Full rewrite of the user manual to document the new mineXpert2 software program. The new mineXpert2 version is a full rewrite of the previous version such that it can handle MS <sup>n</sup> mass spectrometric data. Major features were included in the new version that are described in detail in this user manual.		
Revision 5.8.2	february 2019	fr
Rework the IsoSpec-based isotopic clust calculations to better describe the three available processes. Also, describe the new FWHM calculation feature in the PeakShaperDlg.		
Revision 5.8.0	end january 2019	fr
Add section to described the newly developed feature about isotopic cluster calculations based on the IsoSpec library;		
Revision 5.7.0	beginning january 2019	fr
Finished porting of the historical LaTeX-based documentation to the DocBook/DAPS/FOP publishing system. Please, see Colophon for details;		
Revision not_set	september 2018	fr
Update the document to include new features and to make fixes. Insist on the fact that mineXpert reads mzML data files.		
Revision not_set	march 2018	fr
Update the document to include new m/z~integration features and scripting capabilities.		
Revision not_set	december 2016	fr
Major rewriting of the document to incorporate all the new features. Large chapter on the scripting of mineXpert.		
Revision not_set	december 2017	fr
Refactored document. First version documenting almost all the of current features of the software program.		
Revision not_set	may 2017	fr
Refactored document. First version documenting almost all the of current features of the mineXpert software program.		
Revision not_set	november 2016	fr
Resume writing, with new program name: mineXpert.		
Revision not_set	september 2016	fr
Start of writing.		

# DEDICATION

To Maria Cecilia

To all the admirable people acting in the “*Free Software Movement*” for a better and more ethical computing world

To all involved in the development of the K Desktop Environment (KDE)

To all the readers who helped me with this manual.



# CONTENTS

## PREFACE iv

## I GENERALITIES I

- 1.1 GENERAL CONCEPTS AND TERMINOLOGIES I
- 1.2 ACQUIRING MASS DATA ALONG TIME: TO PROFILE OR NOT TO PROFILE? I
- 1.3 MASS DATA VISUALISATION: TO COMBINE OR NOT TO COMBINE? 2
- 1.4 EXAMPLES OF VARIOUS MASS SPECTRAL DATA INTEGRATIONS 4
  - TIC-> MZ INTEGRATION 4 • TIC-> DT INTEGRATION 5

## 2 THE PROGRAM'S GRAPHICAL USER INTERFACE 6

- 2.1 OPENING MASS SPECTROMETRY DATA FILES 6
- 2.2 THE WINDOW LAYOUT 6
- 2.3 THE MAIN PROGRAM WINDOW MENU 8
- 2.4 THE LOADED MS RUN DATA SETS 8
- 2.5 MASS SPECTRAL DATA LISTING IN A TABLE VIEW 9
- 2.6 THE MAIN DATA-PLOTTING WINDOWS 13
  - THE TIC CHROMATOGRAM WINDOW 15 • THE MASS SPECTRUM WINDOW 16 • THE DRIFT SPECTRUM WINDOW 17 • THE RETENTION TIME VS MASS SPECTRUM COLOR MAP WINDOW 18 • THE DRIFT TIME VS MASS SPECTRUM COLOR MAP WINDOW 19 • THE DRIFT TIME VS RETENTION TIME COLOR MAP WINDOW 20
- 2.7 GENERAL STRUCTURE OF THE PLOT WIDGETS 20
  - THE LEFT COLUMN OF BUTTONS CONFIGURES THE RECEPTION OF A NEW PLOT 21 • THE RIGHT COLUMN OF BUTTONS CONFIGURES THE NEW INTEGRATION TO RUN 23 • THE PLOT WIDGET MAIN MENU 24
- 2.8 GENERAL OPERATION OF THE PLOT WIDGETS 26

<b>3</b>	<b>MASS DATA INTEGRATIONS FEATURED BY MINEXPert2</b>	<b>29</b>
3.1	GENERAL BEHAVIOUR OF PLOT WIDGETS	29
	PROCESSING FLOW ENTITIES DOCUMENT ALL INTEGRATIONS	29 • SETTING
	THE MS <sup>N</sup> FRAGMENTATION PARAMETERS	30 • SETTING THE M/Z
	INTEGRATION PARAMETERS	32 • EFFECTS OF THE M/Z INTEGRATION
	PARAMETERS	35 • REMOVING 0-INTENSITY M/Z DATA POINTS IS
	USEFUL	39 • SAVITZKY-GOLAY FILTERING OF ANY KIND OF DATA
		40
3.2	VARIOUS MASS SPECTRAL DATA INTEGRATIONS	41
	CONSIDERATIONS ON THE DIVERSITY OF MASS DATA CONTENTS	43 • STATISTICAL
	ANALYSIS OF MASS DATA	44
3.3	CHAINED INTEGRATIONS	44
<b>4</b>	<b>MASS SPECTRAL DECONVOLUTIONS</b>	<b>47</b>
4.1	DECONVOLUTION BASED ON CHARGE STATE	47
4.2	DECONVOLUTION BASED ON ISOTOPIC CLUSTER PEAKS	49
4.3	READING THE RESOLVING POWER BASED ON MASS SPECTRAL DATA	50
<b>5</b>	<b>ISOTOPIC CLUSTER CALCULATIONS</b>	<b>52</b>
5.1	CALCULATING ISOTOPIC CLUSTERS WITH IsoSPEC	52
	THE IsoSPEC GRAPHICAL USER INTERFACE IN MINEXPert	53
5.2	SHAPING MASS PEAK CENTROIDS INTO WELL-PROFIED PEAKS	61
<b>6</b>	<b>RECORDING DATA MINING DISCOVERIES</b>	<b>66</b>
<b>A</b>	<b>GNU GENERAL PUBLIC LICENSE VERSION 3</b>	<b>70</b>

# PREFACE

## I SOFTWARE FEATURE OFFERINGS AND INTENDED AUDIENCE

This manual is about the msXpertSuite mass spectrometric software suite, a software environment that contains two modules:

- *massXpert* module: Allows users to define brand new polymer chemistries and use these polymer chemistry definitions to model linear polymer sequences. Once modelled, a polymer sequence can undergo chemical reactions (enzymatic or chemical cleavages, gas-phase fragmentations...). The obtained results are a model of what a mass spectrum would look like if the modelled experiment had actually been carried over up to the mass spectrometry analysis;
- *mineXpert2* module: Allows users to load mass spectrometry data from mzML files, visualize and mine them throughout all their MS<sup>n</sup> depth. The mass data visualization starts either at the TIC chromatogram level or the MS data table view and deepens to the MS<sup>n</sup> mass spectra, the drift spectra, the XIC chromatograms, various color maps relating m/z values with either retention times or drift times. A wide array of mass spectral integrations are available to ease the MS<sup>n</sup> data finest scrutiny.

This manual is therefore intended for people willing to learn how to use the comprehensive msXpertSuite software package.

Mass spectrometry has gained popularity across the past twenty years or so. Indeed, developments in polymer mass spectrometry have made this technique appropriate to accurately measure masses of polymers as heavy as many hundreds of kDa, and of any chemical type.

There are a number of utilities—sold by mass spectrometer constructors with their machines, usually as a marketing “plus”—that allow predicting/analyzing mass spectrometric data obtained on polymers. These programs are usually different from a constructor to another. Also, there are as many mass spectrometric data prediction/analysis computer programs as there are different polymer types. You will get a program for oligonucleotides, another one for proteins, maybe there is one program for saccharides, and so on. Thus, the biochemist/massist, for example, who happens to work on different biopolymer types will have to learn to use several different software packages. Also, if the software user does not own a mass spectrometer, chances are he will need to buy all these software packages.

The msXpertSuite mass spectrometric software is designed to provide *free* solutions to all these problems by providing the following features:

- massXpert:
  - Model *ex nihilo* polymer chemistry definitions (in the XpertDef module that is part of the massXpert program);
  - Perform simple yet powerful mass computations to be made in a mass desktop calculator that is both polymer chemistry definition-aware and fully programmable (that's the XpertCalc module also part of the massXpert program);
  - Edit polymer sequences on a polymer chemistry definition-specific basis, along with chemical reaction simulations, finely configured mass spectrometric computations... (all taking place in the XpertEdit module that is the main module of the massXpert program);
  - Customize the way each monomer will show up graphically during the program operation (in the XpertEdit module);
  - Edit polymer sequences with immediate visualization of the mass changes elicited by the editing activity (in the XpertEdit module);
  - Open an unlimited number of polymer sequences at any given time and of any given polymer chemistry definition type (in the XpertEdit module);
- mineXpert:
  - Load mass spectrometry data files in the mzML format, thanks to the excellent libpwiz library of ProteoWizard<sup>1</sup> fame; Mass data file loading can be performed in *full-memory* mode (all the data read from file are stored in memory) for faster operations or in *streamed* mode when loading files larger than the available memory.
  - Display the data in powerful ways in a unified graphical user interface. The interface was designed to integrate all the most useful characteristics of the various proprietary environments known by the author, thanks to the excellent libqcustomplot<sup>2</sup> library;
  - Configure the way mass spectrometry data integrations are performed from a combination standpoint and optionally apply a Savitzky-Golay smoothing;

---

<sup>1</sup> [HTTP://PROTEOWIZARD.SOURCEFORGE.NET/](http://proteowizard.sourceforge.net/).

<sup>2</sup> [HTTP://QCUSTOMPLOT.COM/](http://qcustomplot.com/).

- Configure the level of MS data ( $MS^n$ ) for which the integrations are performed. Optionally filter data by precursor ion  $m/z$  values or precursor spectrum indices;
- Perform data mining by performing data integrations in various ways;
- Ion mobility mass spectrometry data are supported with a  $int=f(m/z,dt)$  color map plot calculation;
- A specific data integration mode allows easy mass spectral intensity determination for any feature (peak in a TIC chromatogram, a mass spectrum, a drift spectrum, any color map region);

## 2 FEEDBACK FROM THE USERS

We are always grateful to any constructive feedback from the users.

The msXpertSuite software team might be contacted *via* the following addresses

msxpertsuite@msxpertsuite.org = general mailing list about msXpertSuite  
 bug-reports@msxpertsuite.org = report bugs found in msXpertSuite software

FIGURE 1: ADDRESSES TO REPORT FEEDBACK TO

## 3 PROJECT HISTORY

This is a brief history of msXpertSuite.

- 1998–2000

The name massXpert comes from a project I started while I was a post-doctoral fellow at the École Polytechnique (Institut Européen de Chimie et Biologie, Université Bordeaux I, Pessac, France);


The massXpert program was published in *Bioinformatics* (RUSCONI, F. AND BELGHAZI, M. DESKTOP PREDICTION/ANALYSIS OF MASS SPECTROMETRIC DATA IN PROTEOMIC PROJECTS BY USING MASSXPERT. *BIOINFORMATICS*, 2002, 644–655 ([HTTPS://ACADEMIC.OUP.COM/BIOINFORMATICS/ARTICLE/18/4/644/243311](https://academic.oup.com/bioinformatics/article/18/4/644/243311)) ).

At that time, MS-Windows was at the Windows NT 4.0 version and the next big release was going to be “you’ll see what you’ll see”: MS-Windows 2000.

When I tried massXpert on that new version (one colleague had it with a new machine), I discovered that my software would not run normally (the editor was broken). The Microsoft technical staff would advise to “buy a new version of the compiler environment and rebuild”. This was a no-go: I did not want to continue paying for using something I had already produced with legitimate software.

- 2001–2006

During fall 1999, I decided that I would stop using Microsoft products for my development. At the beginning of 2000 I started as a CNRS research staff in a new laboratory and decided to start fresh: I switched to GNU/Linux (I never looked back). After some months of learning, I felt mature to start a new development project that would eventually become an official GNU package: GNU polyxmass.


The GNU polyxmass software, much more powerful than what the initial massXpert software used to be, was published in *BMC Bioinformatics* in 2006 (RUSCONI, F. GNU POLYXMASS: A SOFTWARE FRAMEWORK FOR MASS SPECTROMETRIC SIMULATIONS OF LINEAR (BIO-)POLYMERIC ANALYTES. *BMC BIOINFORMATICS*, 225– (HTTPS://BMCBIOINFORMATICS.BIOMEDCENTRAL.COM/ARTICLES/10.1186/1471-2105-7-226) ).

Following that publication I got a lot of feedback (very positive, in a way) along the lines: —“*Hey, your software looks very interesting; only it’s a pity we cannot use it because it runs on GNU/Linux, and we only use MS-Windows and MacOSX!*”.

- 2007–2016


In december 2006, I decided to make a full rewrite of GNU polyxmass. The software of which you are reading the user manual is the result of that rewrite. I decided to “recycle” the massXpert name because this software is written in C++, as was the first massXpert software. Also, because the first MS-Windows-based massXpert project is not developed anymore, taking that name was kind of a “revival” which I enjoyed. However, the toolkit I used this time is not the Microsoft Foundation Classes (first massXpert version) but the Trolltech Qt framework (see the *About Qt* menu in the *Help* menu in massXpert).

Coding with the Qt libraries has one big advantage: it allows the developer to code once and to compile on the three main platforms available today: GNU/Linux, MacOSX, MS-Windows. Another advantage is that the Qt libraries are wonderful software, technically and philosophically (Free Software).

The rewritten software was published in 2009 (RUSCONI, F. MASSXPERT 2: A CROSS-PLATFORM SOFTWARE ENVIRONMENT FOR POLYMER CHEMISTRY MODELLING AND SIMULATION/ANALYSIS OF MASS SPECTROMETRIC DATA. *BIOINFORMATICS*, 2009, 2741–2742 (HTTPS://ACADEMIC.OUP.COM/BIOINFORMATICS/ARTICLE/25/20/2741/194220) ).

- 2016–2019

In 2016, I started a new project about visualization of mass spectrometric data. The project developed pretty quickly, as we needed at the mass spectrometry facility a software that would allow to cope efficiently with ion mobility mass spectrometric experimental data. mineXpert was thus started.


To bundle both massXpert and mineXpert in a single software suite, I bought the msXpertSuite website [HTTP://MSXPERTSUITE.ORG](http://msxpertsuite.org)  and created that new name.


- 2019–

The mineXpert software was published in 2019. The reviewers acknowledged that it was a useful piece of software but complained about its not handling the MS/MS data. Reference to cite: Filippo Rusconi. mineXpert: Biological Mass Spectrometry Data Visualization and Mining with Full JavaScript Ability. *J. Proteome Res.* 2019, 18, 5, 2254–2259 <https://doi.org/10.1021/acs.jproteome.9b00099>.

In response to the reviewers and for my personal project that I had started a little earlier, I decided to fully rewrite mineXpert into mineXpert2. The new version of the software had to support MS<sup>n</sup> data.

## 4 PROGRAM AND DOCUMENTATION AVAILABILITY AND LICENSE

The programs and all the documentation that are shipped along with the msXpertSuite software suite are available at [HTTP://WWW.MSXPERTSUITE.ORG](http://www.msxpertsuite.org) . Most of the time, a new version is published as source, and as binary install packages for MS-Windows (64-bit systems only). No GNU/Linux binary packages are created outside of the autobuilder of the various distributions. As a Debian Developer, the author creates Debian<sup>3</sup> packages that are uploaded on the distribution servers. These packages are available using the system's software management infrastructure (like using the Debian's `apt` command, for example, or the graphical application).

The software and all the documentation are all provided under the Free Software license *GNU General Public License, Version 3, or later, at your option*. For an in-depth study of the *Free Software* philosophy, I kindly urge the reader to visit [HTTP://WWW.GNU.ORG/PHILOSOPHY](http://www.gnu.org/philosophy) .

---

<sup>3</sup> [HTTP://WWW.DEBIAN.ORG/](http://www.debian.org/) 

## I GENERALITIES

In this chapter, I wish to introduce some general concepts around the mineXpert2 program and the way data elements are named in this manual and in the program.

A mass spectrometry experiment generally involves monitoring the  $m/z$  value of analytes injected in the mass spectrometer along a certain time duration. The  $m/z$  value of each detected analyte is recorded along with the corresponding signal intensity  $i$ , so that a mass spectrum is nothing but a series of  $(m/z, i)$  pairs recorded along the acquisition duration. All along the acquisition, the precise moment at which a given analyte is detected (and its  $(m/z, i)$  pair is recorded), is called the retention time of that analyte ( $rt$ ). This retention time is not to be misunderstood as the drift time of that analyte in an ion mobility mass spectrometry experiment.

### I.1 GENERAL CONCEPTS AND TERMINOLOGIES

Most generally, the mass spectrometer acquires an important number of spectra in, say, one second. But all these spectra are *combined* together, and, on the surface, the massist only sees a “slow” acquisition of 1 spectrum per second. This apparent slow acquisition rate is configurable. At the time of writing, generally 1 spectrum per second is recorded on disk. So, say we record mass spectra for 5 minutes, we would have recorded  $(5 \times 60)$  spectra.

### I.2 ACQUIRING MASS DATA ALONG TIME: TO PROFILE OR NOT TO PROFILE?

As a mass spectrometry user, the reader of this manual certainly has used mass spectrometers where mass spectra are acquired and stored in different ways:

- Mass spectra are acquired and summed—the next to the previous—in such a manner that one is left, at the end of the acquisition, with a single spectrum of which the various peak intensities have been increasing all along the acquisition. Indeed, in this mode, each new spectrum is actually “*combined*” to the previously acquired ones. The resulting mass spectrum that is displayed on screen and that gets ultimately stored on



disk is called a *combined spectrum*. This is typically the way MALDI-TOF mass spectrometers are used when manually acquiring data from samples deposited onto sample plates. We refer to this kind of acquisition as an “accumulation” mode acquisition;

- Mass spectra are acquired and stored on disk as a single file containing all the spectra, appended one after the other. There is no combination of the spectra: each time a new spectrum is displayed on screen, that spectrum is appended to the file.<sup>1</sup> This is typically the case when mass spectra are acquired all along a chromatography run and is generally called a “*profile*” mode acquisition. Note that this profile mode acquisition must not be mistaken as the profile mass peak type that negates the centroid mass peak type.

## 1.3 MASS DATA VISUALISATION: TO COMBINE OR NOT TO COMBINE?

In the previous section, we mentioned *spectrum combination* a number of times. What does that mean, that spectra are “combined” together into a single “combined spectrum”? Say we have 200 spectra that need to be combined together into a *single* spectrum that summatively represents the data of these 200 spectra.

First, a new spectrum would be allocated (*result spectrum*), entirely empty at first. Then, the very first spectrum of the 200 spectra is literally copied into that result spectrum. At this point the combination occurs, according to an iterative process that has the following steps:

- Pick the next spectrum of the 200-spectra dataset;
  1. Pick the first (m/z,i) pair of the currently iterated spectrum;
  2. Look up in the *result spectrum* if a m/z value identical to the m/z value of the current (m/z,i) pair is already present;
  3. If the m/z value is found, increment its intensity by the intensity of the (m/z,i) pair;
  4. Else, if the m/z value is not found, add the current (m/z,i) pair to the result spectrum;
  5. Iterate over all the remaining (m/z,i) pairs of the current spectrum and redo these steps.
- Iterate over all the 198 remaining spectra of the dataset and do the steps above for each single iterated spectrum.

At the end of the two nested loops above, the combined spectrum is still a single spectrum that represents—summatively—all the 200 spectra. This whole process is very computing-intensive, in particular if:

---

<sup>1</sup> Although there certainly is spectrum combination going on in the guts of the software, because the system actually acquires much more spectra than is visible on screen and each newly displayed spectrum is actually the combination of many spectra acquired under the surface.

- The  $m/z$  range is large: there are lots of points in each spectrum, which means that for each new  $(m/z, i)$  pair we need to iterate in the long list of  $m/z$  values that make the result spectrum;
- The resolving power of the mass spectrometer is high: there are many points per  $m/z$  range unit.

When a profile mode acquisition is performed, the user gets an innumerable number of distinct spectra, all appended to a single file. These unitary spectra are virtually unusable if an initial processing is not performed. This initial processing of the spectra is called “*total ion current chromatogram calculation*”. What is it? Let's say that the user has performed a profile mode mass spectrometry acquisition on the eluate of a chromatography column. Now, imagine that the spectrometer stores the mass data at a rate of one spectrum per second and that the chromatography gradient develops over 45 min: there would be a total of  $(45 * 60)$  spectra in that file. The question is: — “*How can we provide the user with a data representation that might be both meaningful and useful to start mining the data?*” The conventional way of doing so is to load all the mass spectra and compute the “*total ion current chromatogram*” (the TIC chromatogram). The analogy with chromatography is evident: the TIC chromatogram is the same as the UV chromatogram unless optical density is not the physical property that is measured over time; instead, the amount of ions that are detected in the mass spectrometer is measured over time. That amount is actually the sum of the intensities of all the  $(m/z, i)$  pairs detected in each spectrum. When mass data are acquired during a chromatography run, often, the total ion current chromatogram mirrors (mimicks) the UV chromatogram<sup>2</sup>. For each retention time (RT) a TIC value is computed by summing the intensities of all the  $(m/z, i)$  pairs detected at that specific RT.

How is this total ion current chromatogram computed? This is an iterative process: from the first spectrum (retention time value 0 s), to the second spectrum (retention time value 1 s) up to the last spectrum (retention time 45 min), the program computes the sum of the intensities of all the spectrum's  $(m/z, i)$  pairs. That computation ends up with a map that relates each RT value with the corresponding TIC value. The TIC chromatogram is nothing but a plot of the TIC values as a function of RT values. In that sense, it is indeed a chromatogram.

mineXpert2 works exactly in this way. When mass spectrometry data are loaded from a file, the TIC chromatogram is computed and displayed. This TIC chromatogram serves as the basis for the mass data mining, as described in this manual. The TIC chromatogram serves as the basis for spectral combinations that can be performed in various ways, and not all formally *combinations*, which is why I prefer the term “*integrations*”. Some of these integrations are described below:

- Integrating data from the TIC chromatogram to a single mass spectrum;
- Integrating data from the TIC chromatogram to a single drift spectrum;

---

<sup>2</sup> Unless eluted analytes do absorb UV light but do not either desorb/desolvate or ionize, or both.

Note that the reverse actions are possible (and indeed necessary for a thorough data mining): selecting a region of a mass spectrum and asking that the TIC chromatogram be reconstituted from there; or selecting a region of a drift spectrum and asking that the TIC chromatogram be reconstituted from there also. Finally, integrations may, of course, be performed from a mass spectrum to a drift spectrum, and reverse.

## I.4 EXAMPLES OF VARIOUS MASS SPECTRAL DATA INTEGRATIONS

In the sections below, the inner workings of mineXpert2 are described for some exemplary mass data integrations. For example, when doing ion mobility mass spectrometry data mining, it is essential to be able to characterize most finely the drift time of each and any analyte. Since each analyte is actually defined as one or more  $(m/z, i)$  pairs, it is essential to be able to ask questions like the following:

- What is the drift time of the ions below this mass peak?
- What are all the drift times of all the analytes going through the mobility cell for a given retention time range?
- What are all the ions that are responsible for this shoulder in the drift spectrum?

### I.4.1 TIC -> MZ INTEGRATION

What computation does actually mineXpert2 do when a mass spectrum is computed starting from a TIC chromatogram region, say between retention time RT minute 7 and RT minute 8.5?

1. List all the mass spectra that were acquired between RT 7 and RT 8.5. In this *spectral set*, there might be many hundreds of spectra that match this criterion, if we think that, in ion mobility mass spectrometry,  $\approx 200$  spectra are acquired and stored individually every second (I mean it, every 1 s time lapse);
2. Allocate a new empty spectrum—the “*combined spectrum*”—and copy into it without modification the first spectrum of the spectral set;
3. Go to the next spectrum of the spectral set and iterate into each  $(m/z, i)$  pair:
  - Check if the  $m/z$  value of the iterated pair is already present in the combined spectrum. If so, increment the combined spectrum’s  $(m/z, i)$  pair’s intensity value by the intensity of the iterated  $(m/z, i)$  pair’s intensity. If not, simply copy the iterated  $(m/z, i)$  pair in the combined spectrum;
  - Iterate over all the remaining  $(m/z, i)$  pairs and perform the same action as above.
4. Iterate over all the remaining spectra of the spectral set and perform step number 3.

mineXpert2 then displays the combined spectrum.

#### 1.4.2 TIC -> DT INTEGRATION

What computation does mineXpert2 actually do when a drift spectrum is computed starting from a given TIC chromatogram region, say between retention time RT minute 7 and RT minute 8.5?

What is a drift spectrum? A drift spectrum (mobilogram) is a plot where the cumulated ion current of the detected ions is plotted against the drift time at which they were detected. Let's see how that computation is handled in mineXpert2:

1. Create a map to store all the (drift time, intensity) pairs that are to be computed below, the (dt,i) map;
2. List all the mass spectra that were acquired between RT 7 and RT 8.5. The obtained list of mass spectra is called the *"spectral set"*;
3. Go to the first spectrum of the spectral set and compute its TIC value (sum of all the intensities of all the (m/z,i) pairs of that spectrum). Get the drift time value at which this mass spectrum was acquired. We thus have a value pair: (dt, i), that is, for drift time *dt*, the intensity of the total ion current is *i*;

At this point, we need to do a short digression: we saw earlier that, at the time of this writing, one of the commercial instruments on which the author of these lines does his experiments stores 200 spectra each second. These 200 spectra actually correspond to the way the drift cycle is divided into 200 bin (time bins). That means that in the retention time range [7–8.5], there are (1.5\*60) complete drift cycles. And thus there are (1.5\*60) spectra with drift time *x*, the same amount of spectra with drift time *y*, and so on for the remaining 198 time bins. Of course, a large number of these spectra might be almost empty, but these spectra are there and we need to cope with them.

The paragraph above must thus lead to one interrogation about the current (dt,i) pair: — *"Has the current dt value be seen before, during the previous iterations in this loop?"*. If not, then create the (dt, i) pair and add it to the (dt,i) map; if yes, get the *dt* element in the map and increment its intensity value by the TIC value computed above;

4. Iterate over all the remaining spectra of the spectral set and perform step number 3.

At the end of the loop above, we get a map in which each item relates a given drift time with a TIC value. This can be understood this way: — *"For each drift time value, what is the accumulated ion current of all the ions having that specific drift time?"*.

At this point, mineXpert2 displays the drift spectrum (mobilogram).

## 2 THE PROGRAM'S GRAPHICAL USER INTERFACE

Data mining, in mass spectrometry, entails, for a large part, the relentless scrutiny of the mass spectra by an expert eye. Without a powerful mass spectrum viewer, capable of numerous data display modes, the expert eye remains powerless. In this chapter, the graphical user interface is described in detail, starting from the opening of a mass spectrometry data file, to the navigation through all the data depth within a large set of graphical representations of the data.

### 2.1 OPENING MASS SPECTROMETRY DATA FILES

To start a session, open one or more mass spectra using the menu *File > Open mass spectrum file(s) in streamed mode* menu. *mineXpert2* understands the mzML, mzXML and MGF formats. The file loading procedures, for these formats, are delegated to the excellent [libpwnz](http://libpwnz) library from the ProteoWizard project.<sup>1</sup> Simple txt, asc data where m/z and intensity values are separated by any character that is neither a newline nor a dot nor a digit (loading is handled by a private parser) can be loaded either from file or directly from the clipboard.

There are two variants of the mass spectrometry file opening menu, one for which the mass data are read from file and stored *totally* in memory and one for which the mass data are read from file in *streamed mode*, used to compute the TIC chromatogram and discarded. The latter mode is useful when the mass data are so large that they cannot fit in memory. Optionally, a table view can be displayed with all the data in the MS run data set that was loaded from disk. The TIC chromatogram or the table view can then be used to access the mass data in the file according to criteria set by the user (retention time range, for example).

### 2.2 THE WINDOW LAYOUT

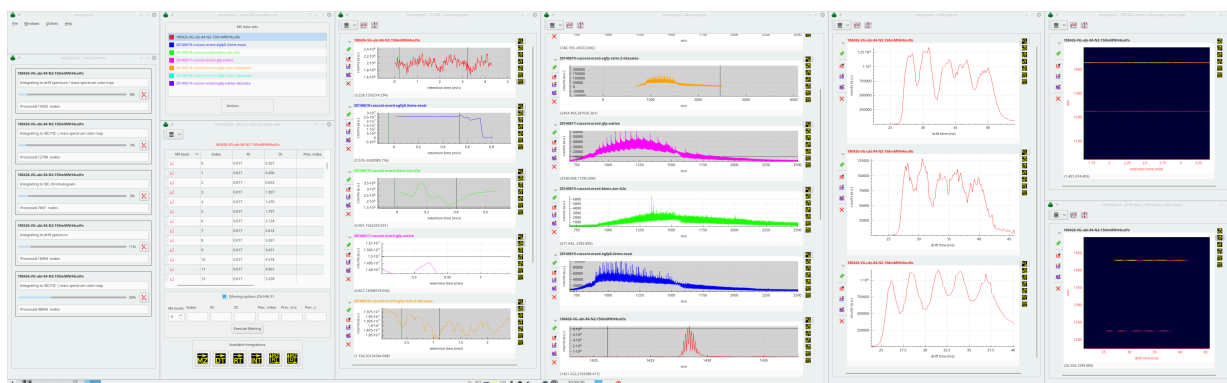
The graphical interface of comprises a number of windows where data and informations are displayed. When first started, the program shows the main program window. Upon loading of a MS run data set from file, new windows do show up as described below (see [FIGURE 2.1, "GENERAL VIEW OF THE GRAPHICAL USER INTERFACE"](#)):

- *mineXpert2* main program window: this is an unintrusive window sporting the main menu;
- The *Task monitors* window, that displays all the ongoing tasks, always providing the user the ability to cancel an ongoing task;
- The *Open MS run data sets* window, that lists all the mass spectrometry data sets that have been loaded from file;

---

<sup>1</sup> Please, see [HTTP://PROTEOWIZARD.SOURCEFORGE.NET/](http://proteowizard.sourceforge.net/) 

- The *TIC/XIC chromatograms* window<sup>2</sup> where the various TIC/XIC chromatograms are displayed for the various mass spectrometry data files that have been loaded. There is, by definition, a single TIC chromatogram for each MS run data set currently loaded in the program. However, this window will also display ion current chromatograms that are computed as a result of an integration step from the other windows, like from the *Mass spectra* window or from the *Drift spectra* window. In this case, the chromatogram is an extracted ion current chromatogram (XIC chromatogram);
- The *Mass spectra* window, where the various mass spectra are displayed. A given mass spectrum may originate from a TIC chromatogram, from the MS run data set tableview window, from a drift spectrum, or even from a color map;
- The *Drift spectra* window, where the various drift spectra are displayed. Like stated before, drift spectra can originate from a variety of places;
- The *TIC/XIC chrom. /Mass spec. color maps* window, that contains  $\text{int} = f(\text{rt}, m/z)$  color maps, two-dimensional representations of the relation between the retention time and the mass spectral data. The intensities are represented as colors in a gradient map;
- The *Drift spec. /Mass spec. color maps* window, that contains  $\text{int} = f(\text{dt}, m/z)$  color maps, two-dimensional representations of the relation between the drift time and the mass spectral data (in ion mobility mass spectrometry experiments only). The intensities are represented as colors in a gradient map;
- The *Console* window, where the various messages or analysis data elements are displayed for the user to select, copy and paste in an electronic lab-book;



The position and size of all the windows can be stored for a later session

**FIGURE 2.1: GENERAL VIEW OF THE GRAPHICAL USER INTERFACE**

<sup>2</sup> TIC stands for “total ion current” and XIC stands for “extracted ion current”.

## 2.3 THE MAIN PROGRAM WINDOW MENU

The menu bar in the main program window displays a number of menu items, reviewed below:

- *File*
  - *File > Open mass spectrum file(s) fully in memory*: Select the mass spectrum file(s) to load. The mass spectral data read from disk are stored in memory. This means that if the file is greater than the available system memory (RAM), the user is urged to select the next menu;
  - *File > Open mass spectrum file(s) in streamed mode*: Select the mass spectrum file(s) to load. The mass spectral data are read from disk, the TIC chromatogram and the MS run data set statistics are computed and then the data are freed. This file opening mode is compulsory when the mass data file is larger than the available system memory (RAM);
  - *File > Load mass spectrum from clipboard*: Create a mass spectrum from a textual representation of (m/z,i) pairs in the same format as described above for the txt,asc file format;
  - *File > Analysis preferences*: Set the preferences for the mining results storage;

- *Windows*

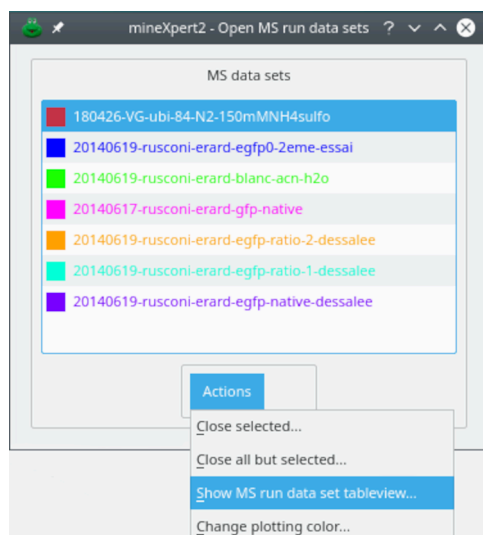
The menus are self-explanatory, as they explicitly explain which window is to be shown. The *Save workspace* menu records on disk the position and size of all the windows, so that upon reopening the program, the windows all position themselves at the recorded position and size;

- *Help*

This menu's items show help about the program itself and also about the Qt libraries that were used to build it. These informations are essential in case the user wants to make a bug report.

## 2.4 THE LOADED MS RUN DATA SETS

Each time a new MS run data set is loaded from disk or from clipboard, a corresponding *MS run data set* item is added to the *Open MS run data sets* dialog window (see FIGURE 2.2, “THE OPEN MS RUN DATA SETS”).



Each time a new MS run data set is loaded, a corresponding item is added to the list widget. The drop-down menu lists a number of actions that can be performed with these items.

**FIGURE 2.2: THE OPEN MS RUN DATA SETS**

One of the most useful menu items available from the drop-down menu of **FIGURE 2.2, “THE OPEN MS RUN DATA SETS”** is the *Show MS run data set table view* item. That menu item triggers the display of a window containing a table view, like in a spreadsheet, that shows the data of each scan contained in the MS run data set loaded from file. This window is described in the following section (see **FIGURE 2.3, “THE TABLE VIEW-BASED SCRUTINY OF ONE MS RUN DATA SET”**).

## 2.5 MASS SPECTRAL DATA LISTING IN A TABLE VIEW

From the *Open MS run data sets* dialog window, it is possible to select one list widget item corresponding to the data set of interest and select the *Show MS run data set table view* menu item that opens the *MS run data set table view* window shown in **FIGURE 2.3, “THE TABLE VIEW-BASED SCRUTINY OF ONE MS RUN DATA SET”**.



mineXpert2 - MS run data set table view

Leptocheline\_MS3\_DDA\_1

MS level	Index	Rt	Dt	Prec. index	Prec. m/z	Prec. z
	104	0.58475667	-1			
	105	0.58859833	-1	104	1521.963013	1
	106	0.59474500	-1	105	1289.949707;1521.963013	1;0;
	107	0.60078167	-1	105	1389.954590;1521.963013	1;0;
	108	0.60782500	-1	105	1189.965088;1521.963013	1;0;
	109	0.61381000	-1	105	1077.947021;1521.963013	1;0;
	110	0.61981000	-1	105	1177.945068;1521.963013	1;0;
	111	0.62525833	-1			

☒ Filtering options

MS levels:  Index:  Rt:  Dt:  Prec. index:  Prec. m/z:  Prec. z:

m/z tolerance for the matches:

Available integrations

Each MS run data set might be scrutinized through the use of this table view. Each filtering criterion has a number of data insertion modalities, described in the text. In this example, the user has set a number of values in the *Filtering options* group box but has not yet executed the filtering, so the full mass spectral data set is listed in the table view.

**FIGURE 2.3: THE TABLE VIEW-BASED SCRUTINY OF ONE MS RUN DATA SET**

The data that belong to each spectrum in the acquired mass spectral data set are shown as values in dedicated columns. Each row represents a mass spectrum of the data set. There might be tens of thousands spectra in an acquisition. Note that two columns may contain -separated values. These columns are *Prec. m/z* and *Prec.z*, respectively precursor ion m/z and precursor ion charge. Ions of differing m/z values and of differing charge values might be accumulated in a trap and then fragmented in a single fragmentation step. The mass data file thus contains spectra with a list of selected precursor ions.



## NOTE

When multiple values are listed in precursor  $m/z$  and charge cells, these values behave as if they were alone in the cell. When performing filtering of the data (see below), then the program treats each value of the list as if it were alone in the cell. For example, in [FIGURE 2.4, “MASS SPECTRAL DATA SET FILTERING USING COMBINED CRITERIA”](#) the user filtered the MS data on the basis of a *Prec.  $m/z$*  value of 1289.95. That filter selected the row that lists, as *Prec.  $m/z$* , the value 1289.952393;1521.973633.

In the *MS run data set table view* window, the checkable group box labelled *Filtering options* allows one to open a series of widgets in which to set the values for various filtering criteria.

- *MS levels*: Set 0 to retain all the mass spectra, irrespective of their MS level. Enter value 1 to filter out any MS level that is not 1. This field has no practical limit, allowing one to perform  $MS^n$  data mining.
- *Index*: Filter by the index number of the mass spectra. This may be of interest when the indices are somehow known in advance.
- *Rt*: Filter by retention time, that is, the time at which the mass spectrum was acquired, relative to the start of the mass data acquisition.
- *Dt*: Filter by drift time (ion mobility mass spectrometry generates this kind of data).
- *Prec. index*: Precursor index.
- *Prec.  $m/z$* : Precursor charge.


The general syntax used to fill-in the filtering options is according to the following scheme:

- *simple value*, like 1 for *MS levels*.
- *<value*: keep rows where values are less than the entered value.
- *>value*: keep rows where values are greater than the entered value.
- *>start <end*: keep rows where values are in between *start* and *end*. The match is not inclusive.
- *start-end*: like before, but the match is inclusive.
- *value%tolerance*: equivalent to the preceding item, unless *start* is (*value* - *tolerance*) and *end* is (*value* + *tolerance*).



## WARNING: THE SPECIFIC CASE OF PRECURSOR M/Z VALUES

When inserting a m/z value, there is a specific tolerance widget *m/z tolerance for the matches* that allows one to define the tolerance with which the matches for the precursor m/z values are performed. That mechanism does not override the *value%tolerance* value definition modality described above; instead, it will add another tolerance to the range values.

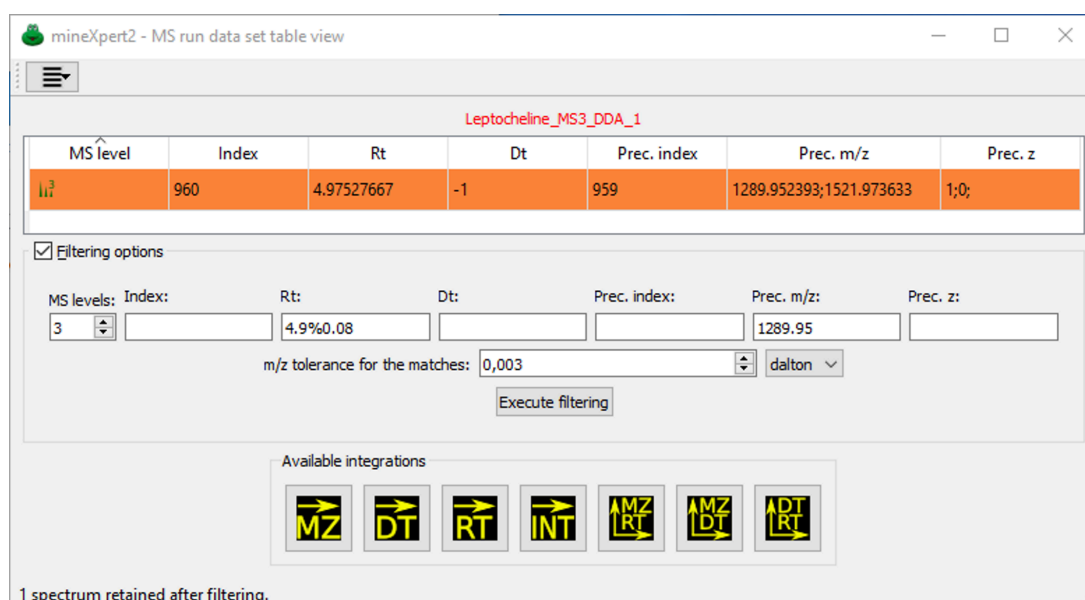
Once the filtering values have been entered and validated by clicking *Execute filtering* or by keying-in , the items (that is, the table rows) that verified the filter, still present in the table view, might be selected (all of them or only some of them) for their integration to a variety of destinations. The destination of each integration is selected by clicking the proper icon push button from the button row below. The meaning of these push buttons will be described later (see [SECTION 2.7.2, “THE RIGHT COLUMN OF BUTTONS CONFIGURES THE NEW INTEGRATION TO RUN”](#)).



## TIP

Make sure to set the m/z integration parameters before performing an integration to a mass spectrum. The setting dialog window opens up using the *Open m/z integration parameters dialog* menu item from the window's main menu, at the top left corner of the window. If the bin size value is too small, the integration might last very much longer than for a “reasonable” size.

The user might enter as many filtering criteria as necessary to pin point the most elusive feature in the mass spectral data set. In [FIGURE 2.4, “MASS SPECTRAL DATA SET FILTERING USING COMBINED CRITERIA”](#), the user has filtered the data to retain MS<sup>3</sup> mass spectra obtained by fragmenting ions of a given m/z value using a specified tolerance for the m/z match in Dalton units.



The use of combined filtering criteria allows one to pin point elusive mass spectral features.

**FIGURE 2.4: MASS SPECTRAL DATA SET FILTERING USING COMBINED CRITERIA**

## 2.6 THE MAIN DATA-PLOTTING WINDOWS

This section will succinctly describe the main data windows of mineXpert2. Each window will be described in greater detail when the features of the program will be described.



## TIP

Although not visible, the various plot widgets in the various plot windows are glued together with splitters that allow their resizing if the window becomes too crowded, as shown in **FIGURE 2.5, “RESIZING THE PLOT WIDGETS WITH THE SPLITTER”**.



When positioning the mouse cursor between two plot widgets, the cursor switches to the splitter cursor shape (circled in red). Dragging that cursor will widen/shrink the plot widgets. To reset the plot widgets as they were initially created, use the *Reset all the plot widget sizes* from the main window menu shown on the right.

**FIGURE 2.5: RESIZING THE PLOT WIDGETS WITH THE SPLITTER**

As visible in all the figures of the plot widget-containing windows described in the next sections, the general structure of these windows is to have a menu in the toolbar, icon buttons in that same toolbar and then plot widgets stacked below in the order in which they have been created.



## NOTE: THE MULTI-GRAPH PLOT WIDGET HAS GONE

In the first version of mineXpert the plot window was divided in two parts; the upper part had a plot widget in which all the plots were overlaid, the lower part had as many plot widgets as there had been integrations because each plot widget could only have one plot in it.

This is no more true now, because, as will be described later, it is now extremely easy to duplicate any plot into any plot widget so as to overlay any number of plots into any number of plot widgets.

## 2.6.1 THE TIC CHROMATOGRAM WINDOW

Each time a new mass spectrum file is loaded, its corresponding TIC chromatogram is computed and then displayed in a new plot widget in the *TIC/XIC chromatogram window* (FIGURE 2.6, “THE TOTAL ION CURRENT (TIC) CHROMATOGRAM WINDOW”). That window is indeed specialized in the plotting of TIC or XIC chromatograms. Each new TIC chromatogram plot that is generated as a result of the loading of a mass spectrometry data file is plotted using a new color. That color encodes the filiation of the whole set of plots that are generated while performing various integrations starting from that initial TIC chromatogram plot. For example, a red TIC chromatogram plot that serves as the starting point for a mass spectrum integration triggers the creation of a new mass spectrum trace plot that is plotted in red color. For color map windows, where there is not trace plot, the axis and tick labels of the plot are of the same color as that of the initial TIC chromatogram plot.



As many traces as necessary can be shown in the window. Each plot has its own set of crosshair markers.

FIGURE 2.6: THE TOTAL ION CURRENT (TIC) CHROMATOGRAM WINDOW

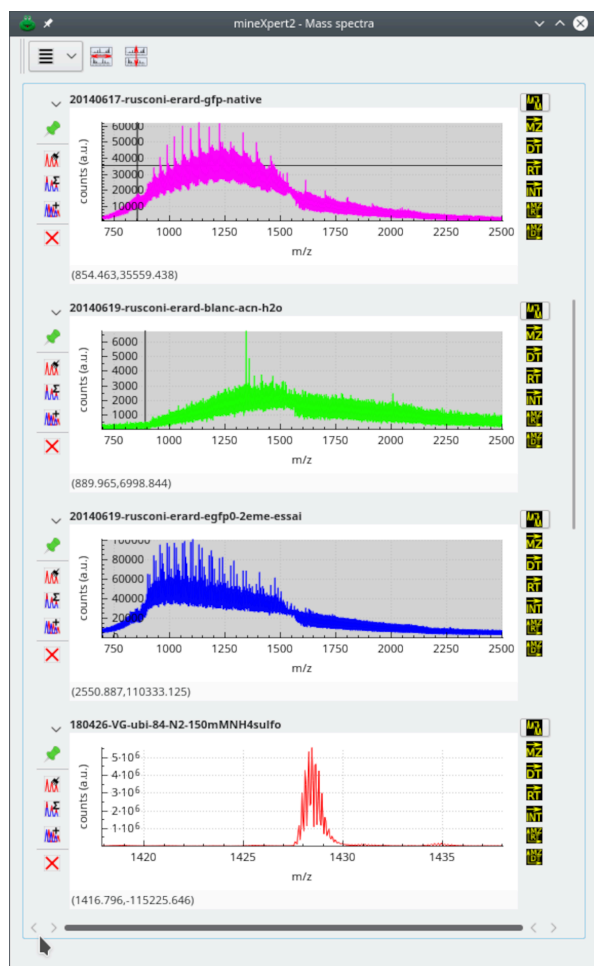


## WARNING: UNCONVENTIONAL SITUATIONS WHILE LOADING MASS SPECTROMETRY DATA FILES

- When the user loads mass spectrometric data from a non-profile acquisition data file or from clipboard data, like when a mass spectrum is opened from a txt,asc,xy text-based format file where the data correspond to a *single* spectrum, *not* a sequence of spectra, the TIC chromatogram really has a single (rt,i) pair denoting the TIC intensity at the single retention time of that very unique spectrum. The TIC chromatogram window thus artificially creates and displays a TIC chromatogram that is a simple line, like shown in the bottom TIC chromatogram plot of [FIGURE 2.6, “THE TOTAL ION CURRENT \(TIC\) CHROMATOGRAM WINDOW”](#). From there, integration happens like in conventional situations. Make sure that the mouse drag movement encompasses the whole TIC region by first unzooming a bit the trace.
- When loading data files that only contain MS<sup>n</sup> data ( $n > 1$ ), like the proteomics-oriented *Mascot generic files* (MGF), mineXpert2 cannot compute a TIC chromatogram. Indeed, a TIC chromatogram can only be computed for MS data, not MS<sup>n</sup> data. In this case, the TIC chromatogram plot that is created does not contain any plot. The only way the user can start mining such kind of MS<sup>n</sup> ( $n > 1$ ) data is by using the MS run data set table view window ([SECTION 2.5, “MASS SPECTRAL DATA LISTING IN A TABLE VIEW”](#)).

### 2.6.2 THE MASS SPECTRUM WINDOW

The mass spectrum window is specialized with the plotting of mass spectra. It thus contains all the plot widgets that display all the mass spectra that were created as a result of mass spectral data integrations. These integration might have originated in any plot widget of any kind (mass spectrum, drift spectrum, color map, for example). For example, the user might select a region in a TIC chromatogram and then ask that a mass spectrum integration be computed. In this case, the resulting mass spectrum is displayed in a new plot widget that is located in the mass spectrum window ([FIGURE 2.7, “THE MASS SPECTRUM WINDOW”](#)).



This mass spectrum window shows multiple mass spectra.

**FIGURE 2.7: THE MASS SPECTRUM WINDOW**

### 2.6.3 THE DRIFT SPECTRUM WINDOW

As described for the mass spectrum window, the drift spectrum window contains all the plot widgets that display drift spectra (FIGURE 2.8, “THE DRIFT SPECTRUM WINDOW”).





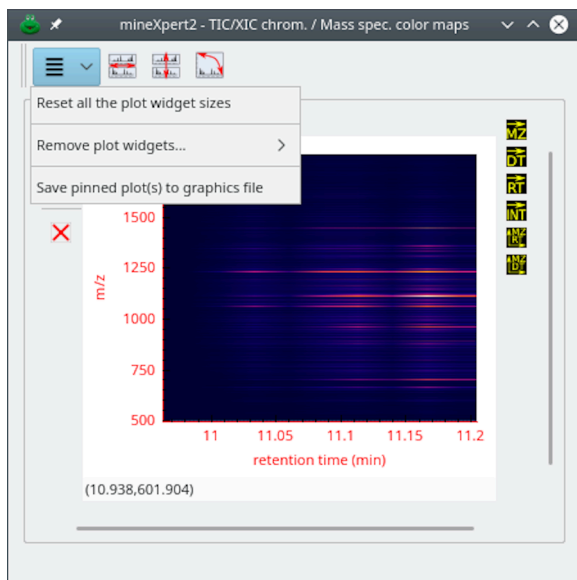
This drift spectrum window shows multiple drift spectra.

**FIGURE 2.8: THE DRIFT SPECTRUM WINDOW**

#### 2.6.4 THE RETENTION TIME vs MASS SPECTRUM COLOR MAP WINDOW

The  $\text{int} = f(\text{rt}, m/z)$  color map window displays a heat map relating all the retention time values with the corresponding mass spectra (see [FIGURE 2.9, “THE  \$\text{INT} = F\(\text{RT}, M/Z\)\$  COLOR MAP WINDOW”](#)). The window main menu offers menus that are all self-explanatory. The tool bar sports three buttons that:

- Lock the x-axis of all the plot widgets in the window (first button).
- Lock the y-axis of all the plot widgets in the window (first button).
- Transpose the axes of the all the plot widgets that are pinned-down (see [NOTE: THE EXPRESSION “PINNED-DOWN WIDGET”](#)).

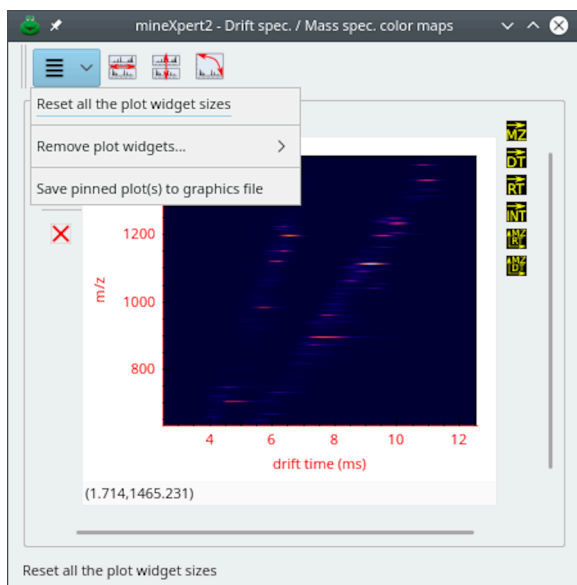


The mass spectrum *vs* retention time color map window relates mass spectra (y-axis) to retention times (x-axis).

**FIGURE 2.9: THE  $\text{INT} = F(\text{RT}, \text{M/Z})$  COLOR MAP WINDOW**

### 2.6.5 THE DRIFT TIME *vs* MASS SPECTRUM COLOR MAP WINDOW

The  $\text{int} = f(\text{dt}, \text{m/z})$  color map window displays a heat map relating all the drift time values with the corresponding mass spectra (see FIGURE 2.10, “THE  $\text{INT} = F(\text{DT}, \text{M/Z})$  COLOR MAP WINDOW”).

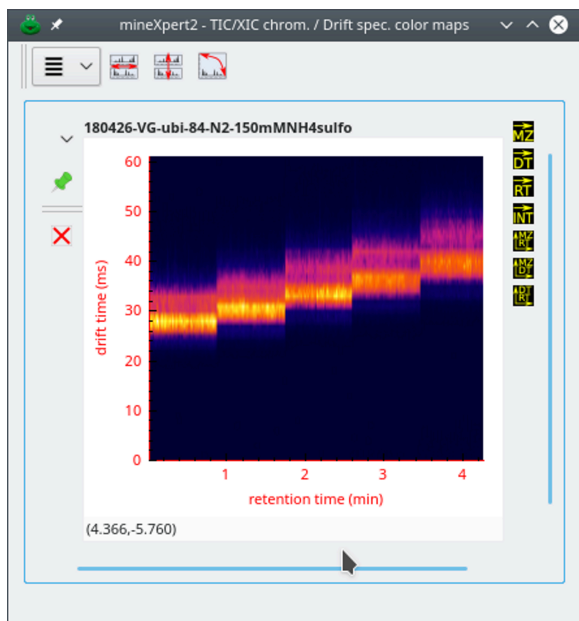


The mass spectrum *vs* drift time color map window relates mass spectra (y-axis) to drift times (x-axis).

**FIGURE 2.10: THE  $\text{INT} = F(\text{DT}, \text{M/Z})$  COLOR MAP WINDOW**

### 2.6.6 THE DRIFT TIME *vs* RETENTION TIME COLOR MAP WINDOW

The  $\text{int} = f(\text{rt}, \text{dt};)$  color map window displays a heat map relating all the drift time values with the corresponding retention time values (see FIGURE 2.11, “THE  $\text{INT} = F(\text{RT}, \text{DT})$  COLOR MAP WINDOW”).

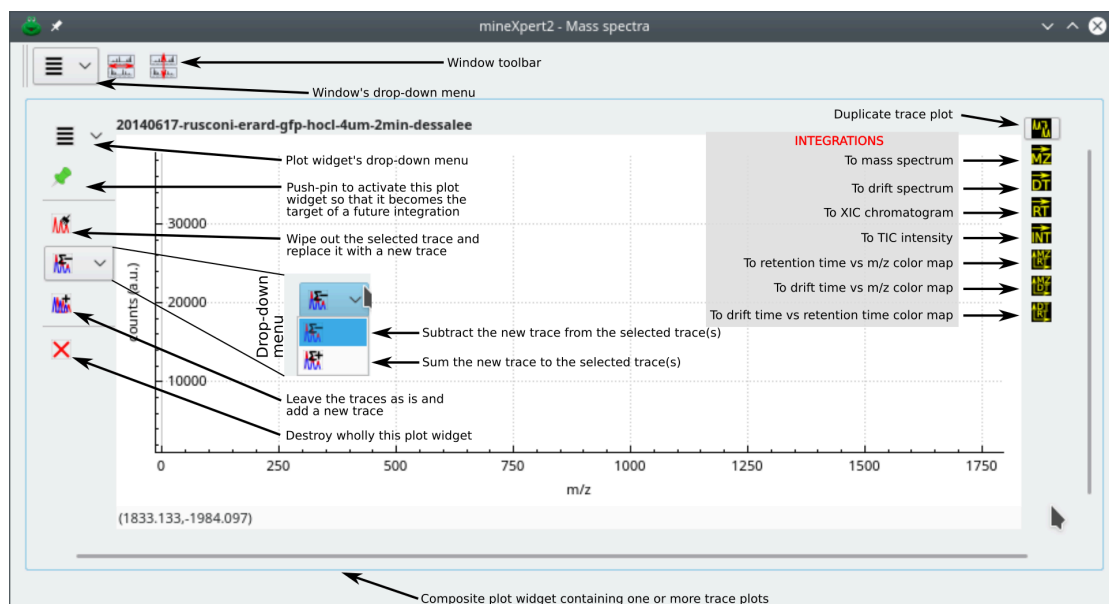


The drift spectrum *vs* TIC chromatogram color map window relates retention times (y-axis) to drift times (x-axis).

FIGURE 2.11: THE  $\text{INT} = F(\text{RT}, \text{DT})$  COLOR MAP WINDOW

## 2.7 GENERAL STRUCTURE OF THE PLOT WIDGETS

The TIC/XIC chromatograms, mass spectra, drift spectra, retention time *vs* mass spectrum and drift time *vs* mass spectrum color map windows (collectively called *plot widget windows*) are all structured in a similar way as described below.





A composite plot widget is a widget that contains a number of other widgets. These sub-widgets are described.

**FIGURE 2.12: DESCRIPTION OF THE VARIOUS WIDGETS THAT MAKE A COMPOSITE PLOT WIDGET**

The window in [FIGURE 2.12](#), “DESCRIPTION OF THE VARIOUS WIDGETS THAT MAKE A COMPOSITE PLOT WIDGET” contains a single composite plot widget. Each composite plot widget's plotting area is sided by two columns of icons, which really are, for most of them, check buttons. The buttons in the left column configure if and how this plot widget is going to be the *receiving container* of a new integration result (result that takes the form of a trace plot). The buttons in the right column configure *what kind of integration* is to be performed using, as a starting point, one of the plots of this plot widget. Indeed, as already mentioned above, integrations in any plot widget might be configured to generate any kind of data: a XIC chromatogram, a mass or drift spectrum... If the data generated by a given integration constitute a XIC chromatogram, that chromatogram will be plotted in the *TIC/XIC chromatograms* window. Both columns of buttons (left and right) are detailed in the next section because these constitute the basis of the operation of mineXpert2.

### 2.7.1 THE LEFT COLUMN OF BUTTONS CONFIGURES THE RECEPTION OF A NEW PLOT


For this description, we will assume that the user is working in the *TIC/XIC chromatograms* window and is performing a new integration to a mass spectrum. The result of that new integration is thus a plot that needs to be installed in the *Mass spectra* window, as depicted in [FIGURE 2.12](#), “DESCRIPTION OF THE VARIOUS WIDGETS THAT MAKE A COMPOSITE PLOT WIDGET”. Such a window, however, might already contain a number of plot widgets; it is thus necessary to establish where the new plot is going to be created and how it will be created. These two questions are answered in detail below by describing the widgets in the left column of widgets in a composite plot widget.

- The plot widget's drop down menu (downward arrow) will be described in detail later. The menu options are actions that apply to the plot widget that owns the menu;
-  The push-pin button, if it is pushed (therefore colored in red ) , indicates that the plot widget “claims” to be the receiving container for any new integration destined to the window containing it.


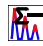




## NOTE: THE EXPRESSION “PINNED-DOWN WIDGET”

A plot widget having its push-pin pushed down is said to be *pinned-down*. This is a terminology that will be used extensively in the rest of this manual.

If the window contains more than one plot widget *and* if more than one plot widget are pinned-down (  ), then all these plot widgets claim to be the receiving container of any new plot to be installed in the window. If no plot widget in a window is pinned-down, any new plot will be created in its own composite plot widget that will appear at the bottom of the plot widget window.








In the list items below, we assume that the plot widget containing the various check buttons is pinned-down and belongs to a non-color map plot window.<sup>3</sup>

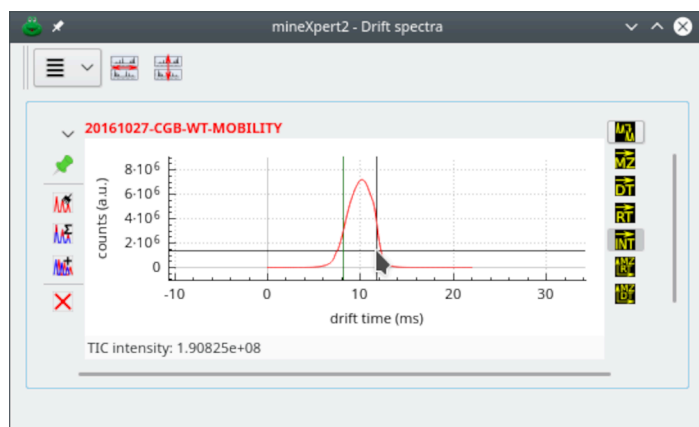
- Checking the  check button indicates that the new plot will *replace* any selected plot in the plot widget.
- Checking the  or the  check button indicates that the new plot will be *subtract-combined* or *add-combined* from or into any selected plot in the plot widget. This feature is particularly useful when trying a step-wise estimation of the mass spectral data out of discrete discontinuous features like TIC chromatogram peaks or drift spectrum peaks. Likewise and reciprocally, it might be useful to extract, from different mass spectral features in a given mass spectrum, the combined TIC chromatogram for them. Typically, when an analyte appears in a mass spectrum under different charge states that are separated by non-baseline regions, being able to assess the XIC chromatogram for them in a series of combined integrations is extremely useful. The subtract-combination is typically used for removing background from a given trace using other blank trace as the background level.
- Checking the  check button indicates that the new plot will be added to the plot widget, that is, it will be overlaid on top of any other plot in the plot widget.

<sup>3</sup> Combination of color maps is not yet implemented, although they might be useful in some situations.

### 2.7.2 THE RIGHT COLUMN OF BUTTONS CONFIGURES THE NEW INTEGRATION TO RUN




Each plot widget may contain zero or more plots (also called graphs or traces). As mentioned earlier, any plot can serve as the starting point for a new integration to the same or another kind of data. For example, from a plot in the *TIC/XIC chromatograms* window, one might want to integrate a given retention time range to a mass spectrum or to a drift spectrum. How are those destinations determined ? The destination of an integration is selected by pressing one of the buttons in the right column of check buttons described below.

-  This button is an exception, because it actually is *not* a check button. It is a true button that, when clicked, elicits the duplication of the currently selected trace(s). There is no integration going on here: the trace(s) are copied “as is” with no processing whatsoever. If none of the plot widgets in the window is pinned-down (that is, all are in the  state ; see **NOTE: THE EXPRESSION “PINNED-DOWN WIDGET”**), the trace(s) are duplicated, each one in its own new plot widget that gets appended to the other widgets at the bottom of the window. If, instead, other plot widget(s) are pinned-down, then the way the duplicated trace(s) will be added depends on the left column buttons. For example, if the  is checked, then the duplicated trace(s) will actually be combined *into* any other selected (or the single present) trace in the pinned-down widget.
-  This is the check button that, when pressed, directs the new integration to produce a mass spectrum to be added to the *Mass spectra* window. Integration to a mass spectrum implies configuring how the m/z values are to be combined (without or with binning and, if so, the size of the bins). Please, refer to **SECTION 3.1.3, “SETTING THE M/Z INTEGRATION PARAMETERS”** for a detailed explanation.
-  This is the check button that, when pressed, directs the new integration to produce a XIC chromatogram to be added to the *TIX/XIC chromatograms* window. “RT” stands for “retention time”.
-  This is the check button that, when pressed, directs the new integration to produce a drift spectrum to be added to the *Drift spectra* window. “DT” stands for “drift time”.
-  This is the check button that, when pressed, directs the new integration to produce a *single TIC intensity value* that is displayed at the bottom of the plot widget, in the widget's status bar.



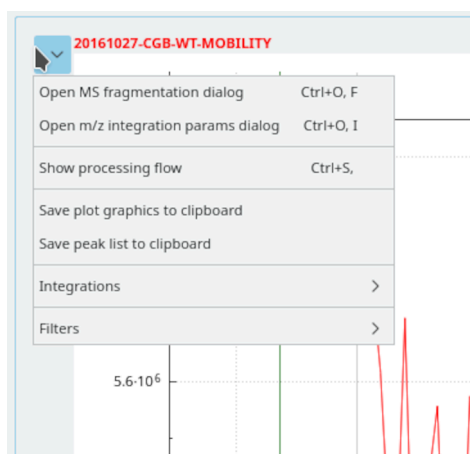
When using the *single TIC intensity value* function, any feature of a plot will be integrated to a single TIC value that represents the sum of all the TIC values for all the MS data in the MS run data set that match the feature selection. The obtained single TIC intensity value is displayed at the bottom of the plot widget in which the calculation was requested. In this example, the feature that was selected for integration is located between the vertical markers.

**FIGURE 2.13: SINGLE TIC INTENSITY VALUE FOR ANY MASS SPECTRAL DATA FEATURE**

-  This is the check button that, when pressed, directs the new integration to produce a  $\text{int} = f(m/z, rt)$  color map. This function is only useful for mass spectrometric data sets that were acquired in “profile” mode; that is, for acquisitions performed along a time frame, like an on-line chromatography-MS experiment.
-  This is the check button that, when pressed, directs the new integration to produce a  $\text{int} = f(m/z, dt)$  color map. This function is only useful when the mass spectra were acquired during an ion mobility mass spectrometry experiment.
-  This is the check button that, when pressed, directs the new integration to produce a  $\text{int} = f(dt, rt)$  color map. This function is only useful for mass spectrometric data sets that were acquired in an ion mobility mass spectrometry experiment.

### 2.7.3 THE PLOT WIDGET MAIN MENU

Each plot widget has a main menu, as depicted in [FIGURE 2.14](#), “**PLOT WIDGET MAIN MENU**”. The various menu items are described below.



The main menu provides actions that may apply only to the plot widget that owns the menu.

**FIGURE 2.14: PLOT WIDGET MAIN MENU**




## TIP

The menu items described below perform actions that apply to individual plots inside the plot widget that owns the menu. mineXpert2 tries to understand the will of the user: when only one plot is in the widget, the menu action applies to that plot, even if it is not currently selected. As soon as there are more than one plot in the widget, the action requires that the target plot be selected. If no plot is selected, a dialog window will instruct the user to select the specific plot of interest.

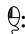

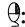
- *Open MS fragmentation params dialog*: Open a dialog to set the MS<sup>n</sup> requirements for the MS data integration (see [SECTION 3.1.2, “SETTING THE MS<sup>N</sup> FRAGMENTATION PARAMETERS”](#)).
- *Open the m/z integration params dialog*: Open a dialog to set the m/z integration parameters (see [SECTION 3.1.3, “SETTING THE M/Z INTEGRATION PARAMETERS”](#)).
- *Show processing flow*: Open a dialog and display the processing flow of the currently selected graph. If there is only one graph in the plot widget, the processing flow for that graph (even if it is not selected) is shown.
- *Save plot graphics to clipboard*: Save the whole plot widget as a graphics object to the clipboard.
- *Save peak list to clipboard*: Save the currently selected (or unique) trace to the clipboard as numerical data.



- *Integrations > Perform single graph point integrations*: integrate the data by moving the cursor over each point of the current trace. Click on the starting trace at the place where point-by-point integration needs to start, then use the  to integrate the next plot point.
- *Filter > Perform Savitzky-Golay smoothing*: the filter will be run with the configuration settings configured as at [SECTION 3.1.6, “SAVITZKY-GOLAY FILTERING OF ANY KIND OF DATA”](#). The smoothed trace will be added into any pinned-down plot widget of the receiving window. If no plot widget is pinned-down, then a new plot widget will be created.


## 2.8 GENERAL OPERATION OF THE PLOT WIDGETS

All the plot widget-containing windows, like the *TIC/XIC chromatograms*, *Mass spectra* and *Drift spectra* windows, all contain plot widgets that have a general working scheme as to how the data can be visualized. The main visualization operations are succinctly described below. The following convention will be used to describe the mouse buttons:

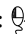

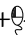
- : left mouse button;
- : middle mouse button;
- : right mouse button.

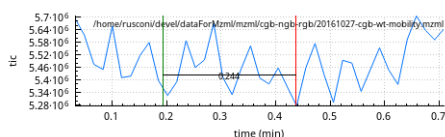


### NOTE

All the visualization operations are performed using the *left* mouse button (). Keyboard keys are used to better qualify the zoom-in/-out or panning operations.

The different plot or colormap graph visualization methods are detailed below:

- Zooming-in and zooming-out:
  - Zoom-in: -click-drag to draw a selection rectangle. When the mouse button is released, the new plot view contains the data contained in the selection rectangle. Note that if the mouse drag over the y-axis spans less than 10 % of the vertical scale, the rectangle does not show up because the software thinks that the mouse drag operation is only horizontal (which has a meaning described later).
  - Zoom-in:  + -click-drag along the x-axis over the region to zoom and release the mouse button. The new zoomed view contains the mouse drag-spanned region. The y-axis spans the region contained between the initial plot bottom y-value and the position on the y-axis of the mouse during the drag operation.



The start and end x locations are visualized with green and red markers, respectively. The x-axis distance between the markers is refreshed all along the mouse movement.

**FIGURE 2.15: ZOOMING-IN OPERATION**

- Zoom-in/-out: **Ctrl** +  $\text{Ⓜ}$ -click-drag *onto* either the x- or the y-axis to interactively zoom-in or -out along that selected axis. In this mode, the zoom operates by contracting/expanding the data in such a manner that the left/bottom part of the graph (the origin of the graph) is anchored and does not move. When the drag occurs towards larger values on the clicked axis, the view is zoomed-in along that axis. Conversely, it is possible to zoom-out by dragging the mouse towards lower axis values.
- Zoom-in/-out: The  $\text{Ⓜ}$ -wheel-rotation can be used to zoom-in or -out the whole plot on both the x- and y-axis simultaneously. Note that the position of the mouse cursor when the wheel is rolled defines the new view of the plot. Practising a bit allows to make that zooming-in/-out mode very powerful.
- Zoom-out: To reset the zoom along one axis,  $\text{Ⓜ}$ -double-click that axis. In this case, only the clicked axis will be full-scale, the other axis remains unchanged. To reset the zoom on both axes in one go,  $\text{Ⓜ}$ -double-click one of the axes maintaining the **Ctrl** key pressed;
- Panning:  
 $\text{Ⓜ}$ -click-drag on one of the axes to pan the plot view along that axis;
- Relative full-scaling along the y-axis: in almost all of the zooming operations described above, it is possible to ensure that the y-axis is full-scaled automatically to the most intense peak visible in the zoomed region by pressing the **Shift** key.
- Recording the history of zoomed views:  
Each time a new zoomed-in/-out view of the plot is triggered, a history element is stored in the plot widget. To back-replay the various steps of the zoom-in/-out operations in sequence, from pre-last to first, hit the **Backspace** key. The exceptions to this mechanics is when the plot view is zoomed using the mouse wheel.



## NOTE: LOCKING THE X AND Y AXES OF ALL THE PLOTS

The tool bar located at the top of all the plot widget-containing windows described above contains two buttons that allow to lock the x-axis (the button icon has the horizontal red line) and the y-axis (red line is vertical) range throughout all the plot widgets in the window. This is of great use when the user wants to compare a number of graphs that have been obtained on comparable samples. The movements and zoom-in or zoom-out operations in one graph are then synchronized to all the other graphs.

### 3 MASS DATA INTEGRATIONS FEATURED BY MINEXPert2

Analyzing mass spectrometric data (with or without drift data) usually involves performing various data integrations in sequence. We saw earlier that upon loading a mass spectrometry data file, the first data visualization that becomes available is the TIC chromatogram. Then, the user might ask to show the MS run data set table view (see [FIGURE 2.3, “THE TABLE VIEW-BASED SCRUTINY OF ONE MS RUN DATA SET”](#)). Thus, the TIC chromatogram and the table view are both starting points for the mass spectrometric data mining. Once a number of mass data integrations have been performed, new integrations will be available in *any* direction, by selecting the integration type in the plot widget (see [SECTION 2.7.2, “THE RIGHT COLUMN OF BUTTONS CONFIGURES THE NEW INTEGRATION TO RUN”](#)).

#### 3.1 GENERAL BEHAVIOUR OF PLOT WIDGETS

There are patterns, in the way mass spectral data integrations can be configured, that are common to all plot widgets. These common patterns are described in the next sections. First, however, the *Processing Flow* concept needs to be described, as it serves as the basis for the integration configurations.

##### 3.1.1 PROCESSING FLOW ENTITIES DOCUMENT ALL INTEGRATIONS

Each time an integration is performed, mineXpert2 stores processing information data that allow to characterize the integration for later internal reuse. For example, any integration documents the kind of computation that needs to be performed. For example, one integration might be from TIC chromatogram to mass spectrum, another integration can be from mass spectrum to drift spectrum. These kinds of integration are documented as processing types. For each integration, there is thus a source and a destination. But this is not sufficient to document precisely the integration: there must be other informational data along with the integration type. When integrating from TIC chromatogram to mass spectrum, the user might select only a region of interest in the chromatogram. That is another kind of information that is documented: the source range (in this case a retention time range). Likewise, when integrating from a mass spectrum to a drift spectrum, the origin is a  $m/z$  range and the destination is a drift spectrum. Other data are stored in the processing information data, like the fragmentation specification—if any has been set—and the  $m/z$  integration parameters, in case the integration produces a mass spectrum.

Together, all the processing information data are stored in what is called a *Processing Flow* entity. A processing flow entity can contain a default  $MS^n$  fragmentation specification and/or default  $m/z$  integration parameters. It can contain any number of *Processing Step* entities that in turn can contain any number of *Processing Spec* entities. Like for the processing flow, each processing step can contain a  $MS^n$  fragmentation specification and/or  $m/z$  integration parameters.

How is the processing flow documented? When a mass spectrometry data file is loaded from disk, mineXpert2 iterates in all the mass spectra of the MS run data set and computes a TIC chromatogram. That TIC chromatogram is displayed in the *TIC/XIC chromatograms* window. At this point, no MS<sup>n</sup> fragmentation specification is set and no m/z integration parameters are set. When the user starts mining the data, their first integration will determine how the processing flow will be filled-in with details about that integration.

If an integration is performed from the TIC chromatogram to a mass spectrum, then m/z integration parameters will necessarily need to be used. Either the user has configured these parameters (SECTION 3.1.3, “SETTING THE M/Z INTEGRATION PARAMETERS”) and they will be used to perform the mass spectral combinations, or default values will be used (these default values are crafted from the statistical analysis of the whole MS run data set). The processing step that is created to document this specific integration will thus document the m/z integration parameters that have been used for the mass spectra combination. That processing step is added to the processing flow that documents the whole integration process. The plot that is created in the *Mass spectra* window will have, associated to it, that processing flow. Think of the *Processing Flow* as a pedigree ID card that each plot has associated to it. That pedigree gets incremented with new processing steps each time an integration is chained to the previous integrations.



## TIP

The original feature of the *Processing Flow* concept is that, because each plot has its own processing flow, when a new integration is performed starting from that plot, then all the preceding steps of the processing flow are replayed and the new step is added to the processing flow. This has the beneficial effect of virtually restricting the scope of the initial data set into a much smaller data set all along the various integrations that occur during a data mining session. When the scope reduces, the computing times reduce accordingly. However, the initial data set is not reduced in memory, only the scope of it is reduced. This is much different from other software in which mass data are effectively pruned off the initial data set with the drawback that data must be reloaded each time an integration must start from the initial data set.

### 3.1.2 SETTING THE MS<sup>N</sup> FRAGMENTATION PARAMETERS

At any moment is it possible to set the MS<sup>n</sup> integration parameters by selecting the *Open MS fragmentation dialog* menu item from the plot widget main menu (SECTION 2.7.3, “THE PLOT WIDGET MAIN MENU”).



The  $MS^n$  fragmentation parameters might be set at any time from the plot widget's main menu. The criteria are of three kinds, as described in detail in the text.

**FIGURE 3.1: SETTING THE  $MS^N$  FRAGMENTATION PARAMETERS**

The dialog window that is displayed is depicted in **FIGURE 3.1**, “SETTING THE  $MS^N$  FRAGMENTATION PARAMETERS”. The available settings are explained below:

- *MS level*: MS level that is targeted by the integration. For example, from a TIC chromatogram, setting *MS level* to 2 would perform an integration only accounting for mass spectra acquired for MS fragmentations of MS level 2, that is, for MS/MS spectra.
- *Precursor m/z values*: `SPACE` or `RETURN` separated list of m/z values. The precursor ions' m/z values are taken into account for the integration. When iterating in the mass spectra acquired in the MS run data set, each spectrum will be checked for the existence of an ion selection list. If any of the m/z values entered in this dialog window are found, the spectrum is accounted for; otherwise it is dismissed. A precursor m/z value is the m/z value of an ion that was selected for fragmentation in a precursor spectrum acquired right before the acquisition of the fragmentation spectrum.



## WARNING

There is a great probability that entering  $m/z$  values in the text widget above will fail to provide any result without some value tolerance set. Indeed, the matches are performed strictly using the entered values and the values recorded by the mass spectrometer, which might have 10 decimals! It is thus necessary to enter a tolerance value for the  $m/z$  value matches. In the figure, the tolerance is set to *0.05 Dalton*. Other “units” are available: resolution power and ppm.

- *Precursor spectra indices*:  or  separated list of integer values. The precursor spectra indices are taken into account for the integration. When iterating in the mass spectra acquired in the MS run data set, each spectrum will be checked for the existence of precursor spectra indices. If the precursor spectra indices entered in this dialog window are found, the spectrum is accounted for; otherwise it is dismissed. A precursor spectrum index is the number of the spectrum in the MS run data set that contained an ion that was selected for further fragmentation.





## NOTE

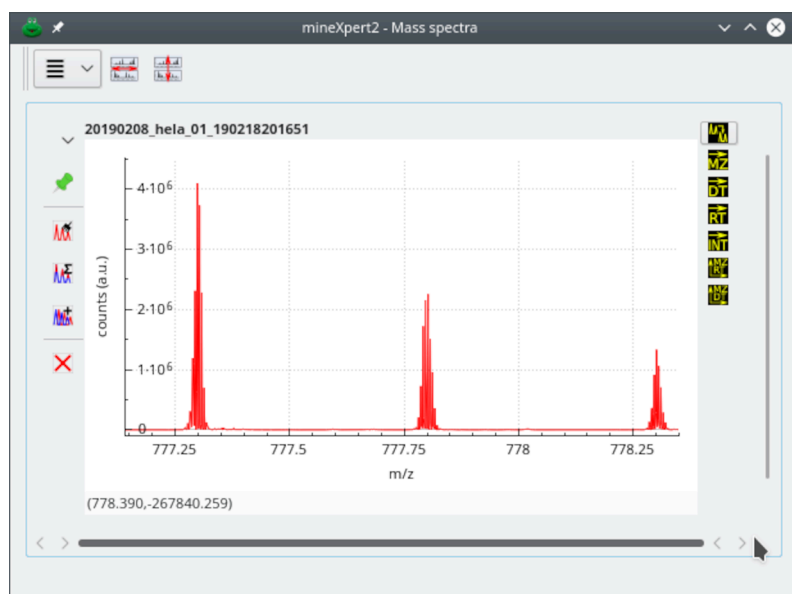
Note that in the dialog window above, none of the three criteria is essential. Setting the *MS level* to 0 deactivates that parameter as a valid criterion for filtering spectra during the integration. In this case, only  $MS^n$  fragmentation specifications used in preceding integrations are taken into account (that is, fragmentation specifications present in the *Processing Flow*). If one wants to reset the  $MS^n$  fragmentation specification to a lower MS level than that of the last preceding integration, the only way to achieve this is to go back to the previous plot where the processing flow contained processing steps of the right  $MS^n$  fragmentation level values. Indeed, it is not possible to reset the  $MS^n$  fragmentation level in a plot by setting the *MS level* to 0 because that plot already has a processing flow containing processing steps configured with a non modifiable MS level.

### 3.1.3 SETTING THE $m/z$ INTEGRATION PARAMETERS

The  $m/z$  integration parameters are only compulsorily used in case an integration will produce a mass spectrum. Indeed, these parameters are required to configure the way mass spectra are combined together into a single product mass spectrum.

An integration to a mass spectrum occurs when the user  selects a range in a given plot while the  checkbox is checked. Integrations to a mass spectrum can be elicited from any plot.

The combination of thousands of spectra to yield a single combined mass spectrum is not something that can happen without tweaking the  $m/z$  integration parameters. Indeed, there is a vast number of different mass spectral data acquisition/storage modes that depend on the vendors or on the mass spectrometer models. In general, this diversity of mass spectral data files<sup>1</sup> creates difficulties in the mass spectra integrations, as shown below in **FIGURE 3.2, “UNUSABLE COMBINATION SPECTRUM WITHOUT BINNING”**.



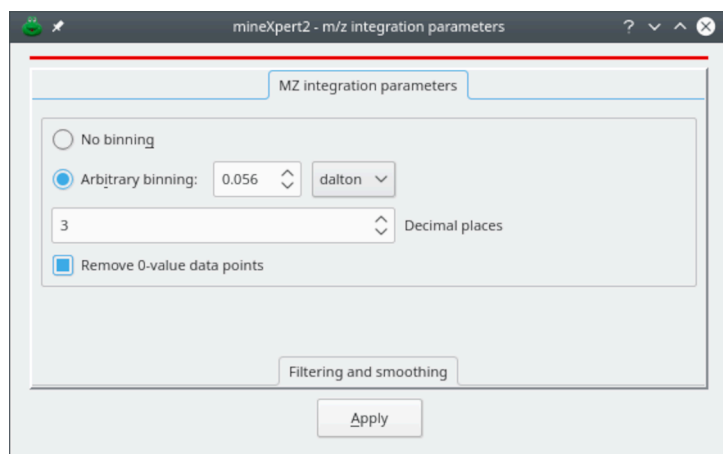
The mass data used to compute this combination spectrum originate from a Lumos Orbitrap analyzer. The visible signal should have been three peaks belonging to a 2-charged ion isotopic cluster. Each peak of the isotopic cluster is artifactually made of numerous data points because the data from the Orbitrap analyzer have 10 decimal digits. Without binning, the mass spectrum is almost unusable.

**FIGURE 3.2: UNUSABLE COMBINATION SPECTRUM WITHOUT BINNING**

In this combination mass spectrum, computed without binning from Lumos Orbitrap-originating data, three peaks that belong to an isotopic cluster appear as made of a number of “sub-peaks”. This is due to the fact that the number of decimals for the  $m/z$  values in the data file is so high that a single isotopic peak appears as a set of peaks. The signal in this mass spectrum is totally useless. This is a perfect illustration of the necessity of data binning of the  $m/z$  values of the mass spectra to be combined.

<sup>1</sup> Even if they are obtained by conversion of raw vendor files to mzML files using ProteoWizard's msconvert tool.





Integrations to a mass spectrum (whatever the source) can be configured to ensure the best results, depending on the kind of mass data. Proper binning configuration is key to getting best results.

**FIGURE 3.3: THE M/Z INTEGRATION PARAMETERS WINDOW**

mineXpert2 provides a number of ways to configure mass spectral combinations such that the obtained mass spectrum is usable. The m/z integration parameters that might be set are described in the following sections. The window depicted in [FIGURE 3.3, “THE M/Z INTEGRATION PARAMETERS WINDOW”](#) can be displayed from any plot widget menu using the *Open the m/z integration params dialog* menu item (see [SECTION 2.7.3, “THE PLOT WIDGET MAIN MENU”](#)). The following parameters can be set:

- *No binning*: no binning is carried-over during the mass spectral integration. In this case, all the m/z values encountered in the mass spectra to be combined are going to be used for the combination and will be encountered in the result spectrum. This is illustrated in [FIGURE 3.2, “UNUSABLE COMBINATION SPECTRUM WITHOUT BINNING”](#).



## TIP

There are instruments that produce perfectly binned mass spectra. In this case, using *No binning* is certainly the best option. For example, the files acquired from the Synapt 2 HDMS mass spectrometer sold by Waters are perfectly binned.

- *Arbitrary binning*: binning is requested and the size of the bins might be defined using either *dalton*, *resolution* or *ppm* “units”.
- *Decimal places*: when performing the combination, each m/z value needs to be restricted to that number of decimal places. This setting may in some situations improve the signal or accelerate the speed of the combination for very large datasets. To leave the decimal places as found in the mass spectrometry data file, enter the value -1.
- *Remove 0-value data points*: remove all the (m/z,i) pairs having a naught *i* value.

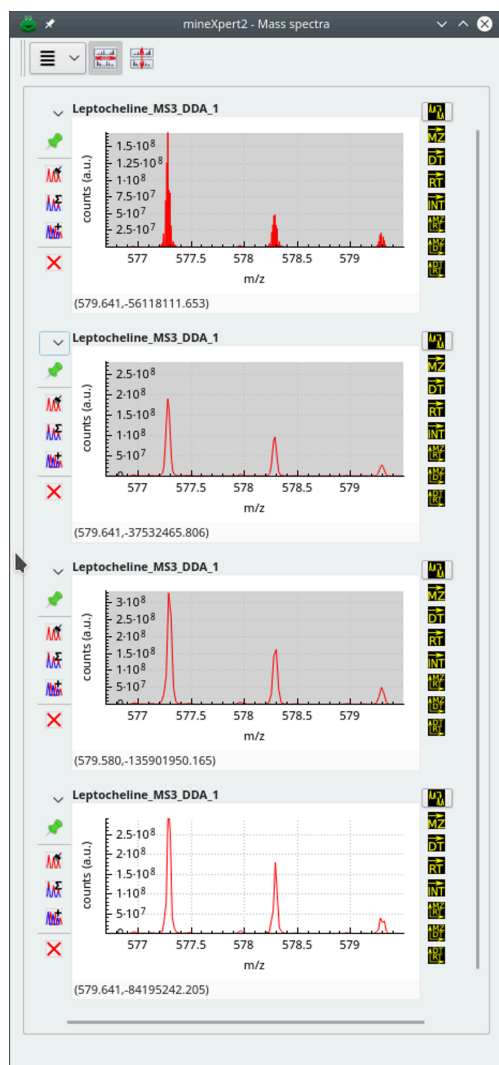
When all the parameters have been set, click onto *Apply* or key in CTRL RETURN . At this point the next integrations started from the plot widget in which these m/z integration parameters were set will be performed using the settings. These parameters will propagate to all the descendant widgets. It will be possible, at any moment and in any widget, to modify the parameters again.

### 3.1.4 EFFECTS OF THE M/Z INTEGRATION PARAMETERS

This section provides some examples of how the m/z integration parameters might impact the mass spectrum resulting from the combination of mass spectra. There might be situations where there is no need to set the m/z integration parameters, typically if the mass data have been properly binned by the software running the mass spectral data acquisition.

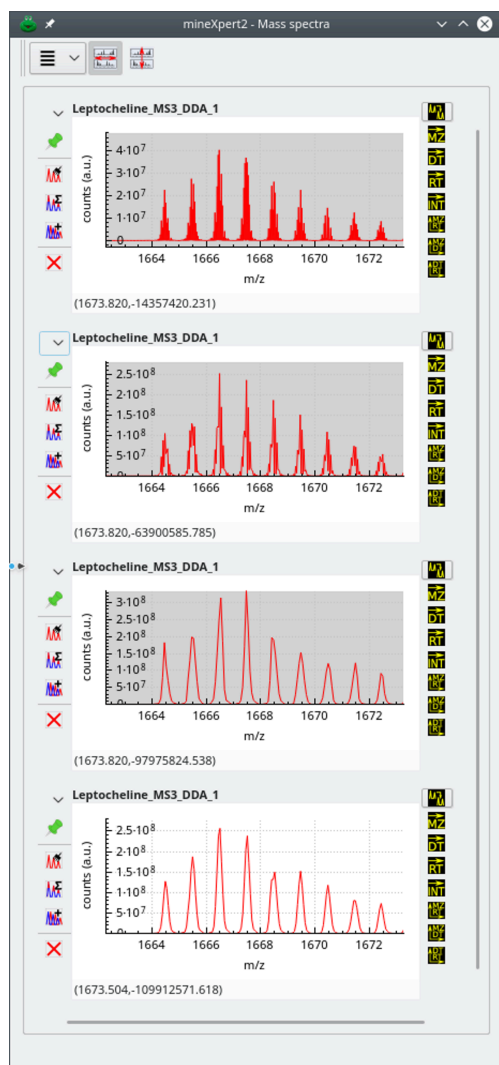
However, in a vast majority of the cases proper setting of the m/z integration parameters will be essential to achieve a combination of mass spectra resulting in a useful mass spectrum. In particular, it is important to grasp that the binning might need to be dynamically set while combining (m/z,i) pairs along the m/z axis. This is why the two sections below show the effects of binning where the bin size unit is either *Dalton* (bins are of fixed-size throughout of the whole m/z range) or either *res* or *ppm* (for resolving power or part-per-million, respectively; bins have varying sizes throughout of the whole m/z range).

In the series of figures below, the same mass spectral data set was integrated to a mass spectrum using different m/z integration parameters. Also, each figure represent a highly zoomed in region of the m/z range. The first figure shows a low m/z region, the second a middle m/z range and finally, the third figure a high m/z range. The figures illustrate the shortcomings of using fixed size bins throughout the whole m/z range of the mass spectrum and the advantage that dynamic ppm- or resolving power-based bin size determination shows.



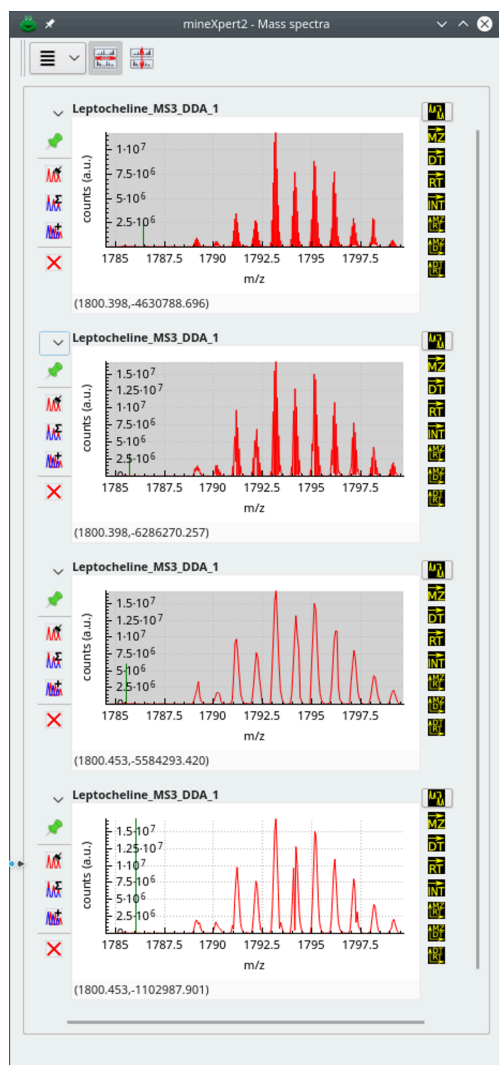
The zoomed in region of the whole m/z range is a low m/z region. The top mass spectrum corresponds to no binning, and below, from top to bottom, binning with sizes: 20 ppm, 40 ppm and 30000 of resolution power. Data from <FTP://MASSIVE.UCSD.EDU/MSV000084765/>.

**FIGURE 3.4: EFFECTS OF BINNING SETTINGS ON THE COMBINATION MASS SPECTRUM (LOW M/Z REGION)**



The zoomed in region of the whole  $m/z$  range is a middle  $m/z$  region. The top mass spectrum corresponds to no binning, and below, from top to bottom, binning with sizes: 20 ppm, 40 ppm and 30000 of resolution power. Data from <FTP://MASSIVE.UCSD.EDU/MSV000084765/>.

**FIGURE 3.5: EFFECTS OF BINNING SETTINGS ON THE COMBINATION MASS SPECTRUM (MIDDLE  $M/Z$  REGION)**



The zoomed in region of the whole  $m/z$  range is a high  $m/z$  region. The top mass spectrum corresponds to no binning, and below, from top to bottom, binning with sizes: 20 ppm, 40 ppm and 30000 of resolution power. Data from <FTP://MASSIVE.UCSD.EDU/MSV000084765/>.

**FIGURE 3.6: EFFECTS OF BINNING SETTINGS ON THE COMBINATION MASS SPECTRUM (HIGH  $M/Z$  REGION)**

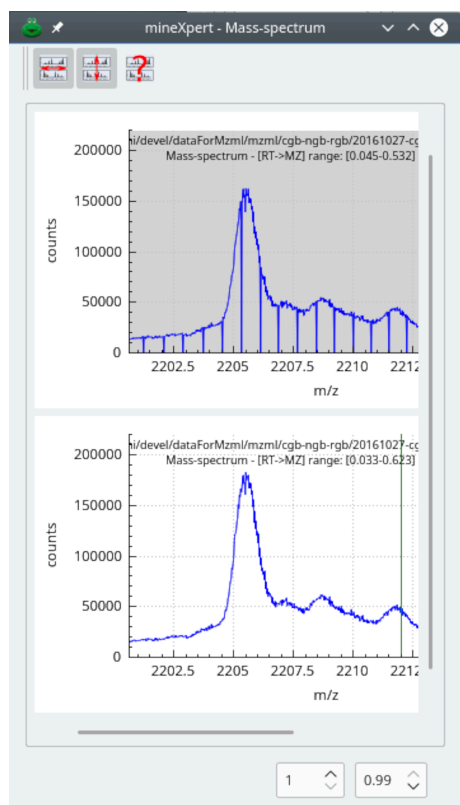
The first observation that can be made by looking at the three figures above is that, for the mass spectral data set at hand, the *No Binning* setting does not provide a usable mass spectrum, whatever the  $m/z$  region of that spectrum's  $m/z$  range. Another interesting observation is that the effect of a given binning setting is not constant throughout the whole  $m/z$  range. For example, if we look at the second spectrum from the top, that was combined using dynamic bins of a 20 ppm width, we see that the spectrum in the low  $m/z$  region is usable, that the spectrum in the middle region is a little less usable and that the spectrum in the high region is not usable at all (too many points artifactually define a mass peak). On the contrary, the spectrum that was generated using larger bins (width of 40 ppm) looks excellent in all three regions of the  $m/z$  range. A similar (although slightly less

good) result was achieved by using also dynamic size bins, but with the resolving power “unit”: 30000 resolving power “units”. Using a lesser resolving power value, like 25000 or 20000, would provide results as good as for the 40 ppm setting above.

### 3.1.5 REMOVING 0-INTENSITY M/Z DATA POINTS IS USEFUL

Mass spectrum data points with an intensity value equal to 0 might arise either from the initial data, as read from the file, or from the binning process, as detailed below.

When bins are prepared, prior to starting a mass spectral combination operation, they all correspond to a m/z value associated to a 0-value intensity. Once the bins have been prepared, the combination starts and, while iterating through all the data points of all the spectra, fills-in the bins. If, however, a bin is never filled (it ends up never updated during the combination), that bin will have an intensity of 0. Bins having a 0-intensity value have a bad effect on the result mass spectrum. Removing them by setting the *Remove 0-value data points* parameter proves beneficial, as shown in **FIGURE 3.7, “REMOVING 0-INTENSITY DATA POINTS”**. When the 0-intensity data points are not removed (upper spectrum), the signal is deteriorated by spurious inverted spikes. Removal of the 0-intensity data points, cleans the trace perfectly.

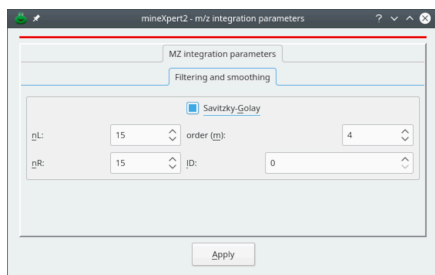


When arbitrary binning is performed, residual 0-intensity data points might survive in the combination spectrum, which deteriorates the resulting mass spectrum. Removing these data points from the combined mass spectrum cleans the trace.

**FIGURE 3.7: REMOVING 0-INTENSITY DATA POINTS**

### 3.1.6 SAVITZKY-GOLAY FILTERING OF ANY KIND OF DATA

The Savitzky-Golay filtering method is widely known for its effectiveness in removing noise from mass spectral data. It is possible to apply that filter to any plot by selecting the corresponding menu as described at [FIGURE 2.14](#), “[PLOT WIDGET MAIN MENU](#)”. The Savitzky-Golay parameters can be set by selecting the *Filtering and smoothing* widget as show on [FIGURE 3.8](#), “[SMOOTHING TRACES WITH THE SAVITZKY-GOLAY FILTER](#)”.



The Savitzky-Golay algorithm is a powerful filter that allows one to smooth any graph. The configuration of the filter is described in the text.

**FIGURE 3.8: SMOOTHING TRACES WITH THE SAVITZKY-GOLAY FILTER**

- $nL$ : specifies the number of data points to the left of the point being filtered;
- $nR$ : specifies the number of data points to the right of the point being filtered. The total number of points in the window that is considered for the regression is thus  $nL + nR + 1$ .
- $m$ : specifies the order of the polynomial to use in the regression analysis leading to the Savitzky-Golay coefficients (typically between 2 and 6);
- $ID$ : specifies the order of the derivative to extract from the Savitzky-Golay smoothing algorithm (for regular smoothing, use 0);


## 3.2 VARIOUS MASS SPECTRAL DATA INTEGRATIONS

Now that the configuration patterns common to all plot widgets have been described, the various mass spectral data integrations can be reviewed.

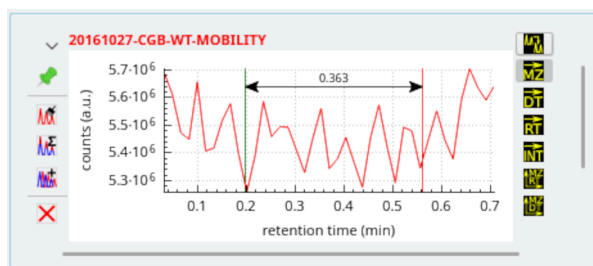


### NOTE

All the integration operations are performed using the *right* mouse button ( $\text{Ⓜ}$ ). This setting allows to easily distinguish all these integration operations from all the  $\text{Ⓛ}$ -based visualization operations (SECTION 2.8, “GENERAL OPERATION OF THE PLOT WIDGETS”).

- *Integrations to a mass spectrum* This kind of operation is triggered upon  $\text{Ⓛ}$ -click-dragging the mouse over the region of interest after having selected the  button on the right button column of the widget. mineXpert2 integrates all the spectra that have been acquired between the start marker and the end marker as set during the mouse drag movement. A new mass spectrum is then plotted in the *Mass spectra* window.



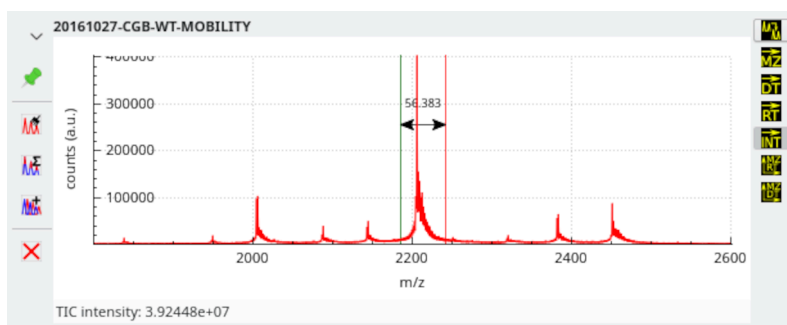


The integration is defined to be a *[TIC chromatogram to mass spectrum]* integration because the plot widget is in the *TIC/XIC chromatograms* window and the *MZ* push button is checked.

**FIGURE 3.9: INTEGRATING MASS DATA FROM A TIC CHROMATOGRAM TO A MASS SPECTRUM**

As seen on the figure above, the region defined by the  $\mathcal{Q}$ -click-dragging operation is delimited by arrows, a green marker at the start and a red marker at the end. The arrows at the ends of the horizontal line allow one to make the difference between an integration selection and a simple distance-measuring selection (detailed later).

- *Integrations to a drift spectrum* This kind of operation is similar to the one described above, unless for the  $\mathcal{DT}$  button that is now checked. The new plot is added to the *Drift spectra* window.
- *Integrations to TIC intensity* This operation involves  $\mathcal{Q}$ -selecting a plot region of interest encompassing the mass spectral feature for which the intensity is to be integrated while the  $\mathcal{INT}$  is checked. The obtained numerical result is displayed in the status bar below the plot widget and in the *Console window*.



The integration is defined to be a *Mass spectrum to TIC intensity* integration because the plot widget is in the *Mass spectra* window and the *INT* push button is checked. The TIC intensity value is printed in the status bar.

**FIGURE 3.10: INTEGRATING MASS DATA FROM A MASS SPECTRUM TO A SINGLE TIC INTENSITY VALUE**

- The same mechanics is at work in the other plot widget windows. For example, to trigger the integration of any kind of plot to a intensity =  $f(dt, m/z)$  color map, simply check the  $\mathcal{MZDT}$  button and drag the mouse over the plot region of interest.

### 3.2.1 CONSIDERATIONS ON THE DIVERSITY OF MASS DATA CONTENTS

Loading data from mass data files in mzML format does not guarantee that the data will be of the same kind when they originate from different mass spectrometers. For example, data from Orbitrap mass spectrometers have the following characteristics:

- All spectra do not start at the same m/z value;
- All spectra do not have the same number of data points (they do not have the same size);
- A large number of data points might have 0 values (intensity at a given m/z value is 0);
- The m/z delta between two consecutive m/z values is not constant, and this is the major difficulty for data integration to a mass spectrum.

This is the output of the statistical analysis of the data loaded from a Lumos Orbitrap-originating file:

Spectral data set statistics:

Total number of spectra: 6203

Average of spectrum size: 391.311946

StdDev of spectrum size: 168.062934

Minimum m/z value: 400.007111

Average of first m/z value: 401.448935

StdDev of first m/z value: 1.590049

Maximum m/z value: 1999.928589

Average of last m/z value: 1901.852315

StdDev of last m/z value: 45.864131

Minimum m/z shift: -0.344452

Maximum m/z shift: 0.000000

Average of m/z shift: 1.097372

StdDev of m/z shift: 1.590049

Smallest Delta of m/z (step): 0.006195

Average of smallest Delta of m/z (step): 0.023757

StdDev of smallest Delta of m/z (step): 0.013179

Greatest Delta of m/z (step): 405.356934

Average of greatest Delta of m/z (step): 163.112057

StdDev of greatest Delta of m/z (step): 75.947334

As mentioned earlier, the most interesting bit of information is in the line reproduced below:

Smallest Delta of m/z (step): 0.006195

That 0.0062 value somehow gives an indication of the “definition” of the spectrum, that is, the smallest distance possible between two consecutive points in the m/z-axis.

In general, the fact that the spectra of an acquisition do not all have the same m/z vector as the m/z-axis is a great difficulty for mass spectral integration because it requires setting up binning prior to performing the mass spectral combination. That binning is nothing else than crafting a m/z value vector able to receive the intensities of all the m/z data points in the spectra to be combined. These concepts are developed in the next paragraph.

### 3.2.2 STATISTICAL ANALYSIS OF MASS DATA

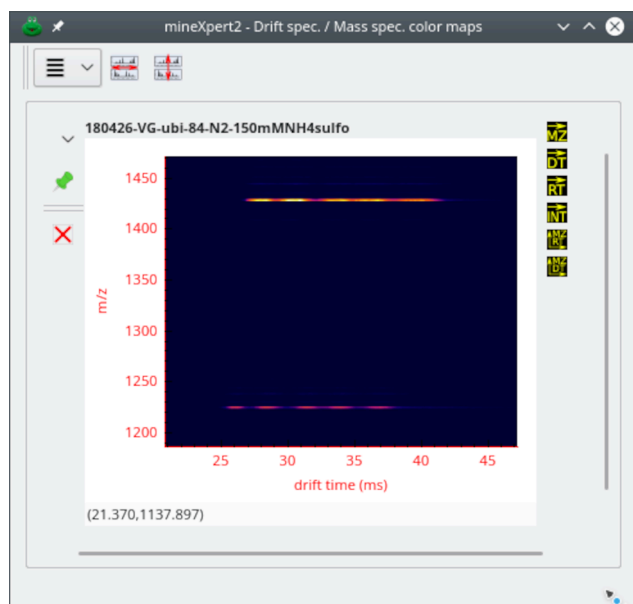
At the end of the data file loading, mineXpert2 performs a rudimentary statistical analysis of the data. The main datum of interest is the smallest m/z step that is observed in the whole set of mass data loaded from disk (the mass spectrum list, that can hold mass spectra in the thousands). For each mass spectrum in the list, the smallest m/z delta between any two consecutive data points is recorded. Then, the smallest ever m/z delta value is sought amidst all the recorded values. Intuitively, that smallest m/z delta value provides an idea of the resolution power of the instrument that generated the mass spectra. Interestingly, this is not the proper value to configure binning. The best value is the median value of the smallest m/z delta values encountered over all the mass spectra of the data file. It is the value that is suggested by default to arbitrarily construct the bins during an integration to a mass spectrum, as described in [FIGURE 3.3, “THE M/Z INTEGRATION PARAMETERS WINDOW”](#) (*Arbitrary binning value with bin size unit dalton*).

## 3.3 CHAINED INTEGRATIONS

The user, in the process of mining the data, will inevitably chain integrations to pinpoint a specific feature of interest. For example, let's say that the user is mining ion mobility mass spectrometry data.<sup>2</sup> After having loaded the data file, the TIC chromatogram is computed and displayed. From there, it is possible to perform a TIC chromatogram to an intensity =  $f(dt, m/z)$  color map integration (see [FIGURE 2.10, “THE INT = F\(DT, M/Z\) COLOR MAP WINDOW”](#).)

---

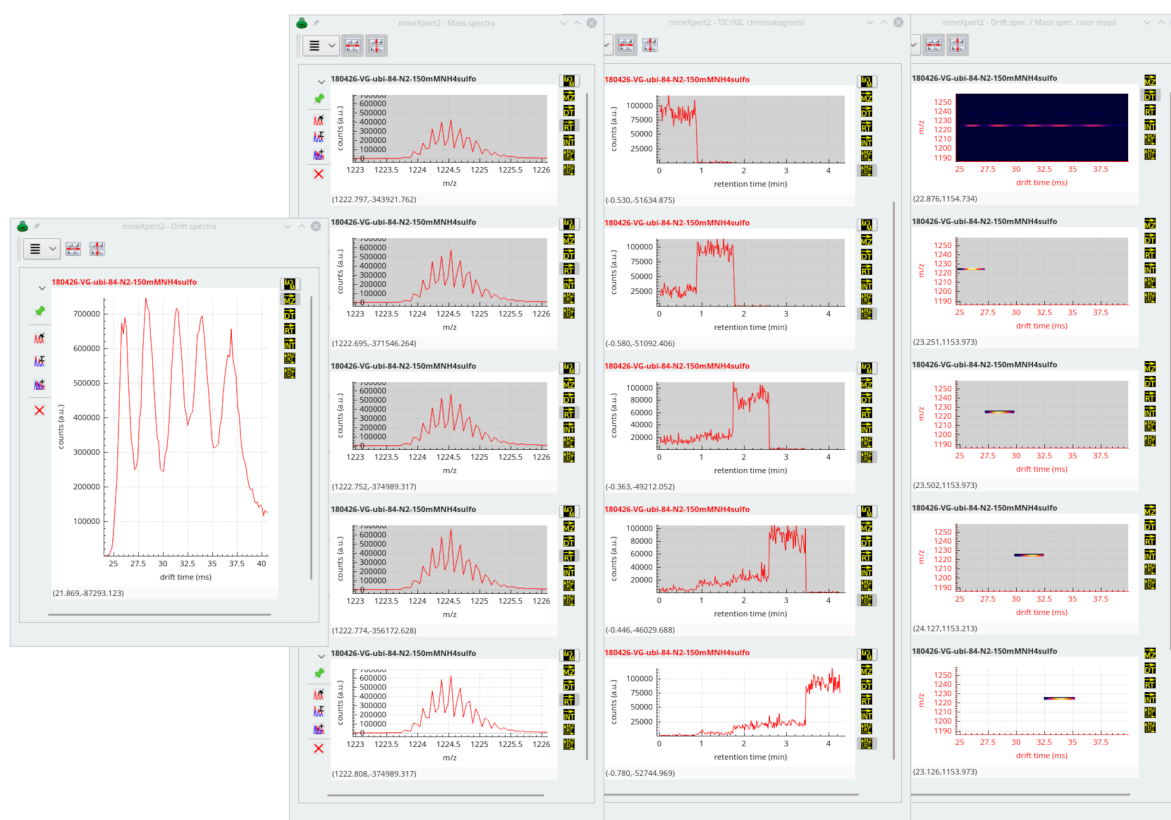
<sup>2</sup> Thanks to Dr Valérie Gabelica (IECB, Bordeaux, France) for permission to use her data set.



**FIGURE 3.11: EXAMPLE OF CHAINED INTEGRATIONS, COLOR MAP WHERE DATA MINING MAY START**

There starts the exploration. The user sees that there are a number of species having discrete drift times at the  $m/z$  ratio around 1220 (lower region of the colormap). They thus integrate to a single drift spectrum that horizontal lower region of the colormap. The obtained drift spectrum is shown at the left hand side of [FIGURE 3.12](#), “**MULTIPLE CHAINED INTEGRATIONS**”.

Because there are five drift peaks in the drift spectrum, the user performs as many individual mass data integrations to a mass spectrum, from left to right. The obtained mass spectra are all shown in the window next to the *Drift spectra* window of the same figure.



**FIGURE 3.12: MULTIPLE CHAINED INTEGRATIONS**

Most interestingly, the various drift regions are integrated to almost identical  $m/z$  values in their respective mass spectrum. In order to know when the various molecular species eluted in the chromatogram, the user performs for each mass spectrum an integration to a XIC chromatogram. The XIC chromatograms are shown on the window next to the *Mass spectra* window. Visibly, each molecular species was eluting from the chromatography column at discrete retention times (this was clearly not a *true* chromatography but instead an infusion during which instrument parameters were changed to modify the mobility properties of the ubiquitin molecule).

At this point, it is interesting to confirm that each XIC chromatogram contains exactly the molecular species that were responsible for the appearance of the purple “bands” on the  $\text{int} = f(\text{dt}, m/z)$  initial color map. The different XIC chromatograms are thus integrated to a color map and the results are shown below the initial color map (less one color map that did not fit in the window). Visibly, the extracted color maps reconstitute the initial pattern visible at the top of the *Drift spec / mass spec color maps* window.



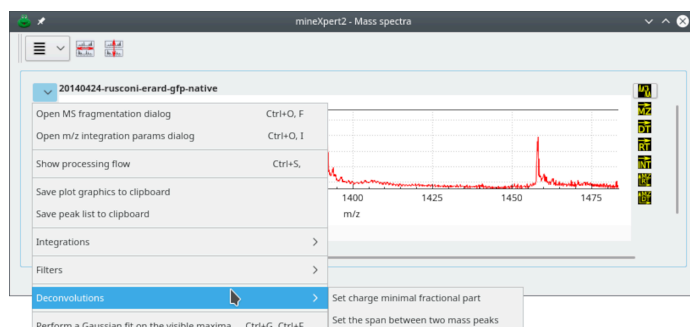
## NOTE

The dataset used for **SECTION 3.3, “CHAINED INTEGRATIONS”** was kind courtesy of Dr. Valérie Gabelica and correspond to a work entitled *Optimizing Native Ion Mobility Q-TOF in Helium and Nitrogen for Very Fragile Noncovalent Structures* published in *JASMS* with DOI: 10.1007/s13361-018-2029-4.

## 4 MASS SPECTRAL DECONVOLUTIONS

When analysing a mass spectrum, two major deconvolutions are performed to get back to the  $M_r$  mass of the analyte while reading  $m/z$  values: the charge-based deconvolution and the monoisotopic cluster-based deconvolution. In the following sections, both deconvolutions are described.

Before delving in the deconvolutions, it is necessary to present two menu options that are found in the plot widgets contained in the *Mass spectra* window: the menu items under the *Deconvolutions* menu (FIGURE 4.1, “MASS SPECTRUM PLOT WIDGET-SPECIFIC DECONVOLUTION MENU”).



The two menu items are needed to configure the mouse-based deconvolution of mass spectra.

FIGURE 4.1: MASS SPECTRUM PLOT WIDGET-SPECIFIC DECONVOLUTION MENU

These two menus allow one to set parameters for the deconvolution (see text for details).

### 4.1 DECONVOLUTION BASED ON CHARGE STATE

In this kind of deconvolution, at the present time, the software assumes that the ionization agent is the proton and that the ionization is positive.

The deconvolution is based on the determination of the distance between two peaks —consecutive or not— of a given charge state envelope. When the user  $\mathcal{O}$ -click-drags the cursor from one peak to another, the program tries to calculate if the distance between two peaks matches one or more charge difference(s). If so, it computes the molecular ( $M_r$ ) mass of the analyte whose mass peak is located *under the cursor*.



#### NOTE: THE MOUSE DRAG POSITION IS SIGNIFICANT

Note that the  $\mathcal{O}$ -click-dragging direction (left→right or right→left) has an impact on the value of the charge  $z$  that is obtained, since that charge value is computed for the peak located *under* the cursor at the moment of the deconvolution. Conversely, the mouse-dragging direction has no effect on the  $M_r$  ( $[^{12}\text{C}]$ -relative molecular mass) of the analyte obtained as a result of the deconvolution process.

FIGURE 4.2, “CHARGE STATE-BASED MASS DECONVOLUTION” shows that process for a protein of  $M_r \approx 30599$  Da. In the top panel, the deconvolution has involved two consecutive peaks. This is the default setting. However, sometimes in low-amount sample mass spectra, two nicely configured consecutive mass peaks are not found, and it is necessary to search for peaks at more than one peak span. For that, use the *Set the span between two mass peaks* to the required value (see FIGURE 4.1, “MASS SPECTRUM PLOT WIDGET-SPECIFIC DECONVOLUTION MENU”). In the lower panel, that span value was set to 2 and the deconvolution provided the same result (of course, since we went one more peak on the left side of the spectrum, that pointed peak corresponds to an ion bearing one more proton).



Deconvolution approach using two peaks belonging to the same charge state envelope. The top deconvolution involves two consecutive mass peaks (peak span value is 1). The bottom deconvolution involves two non-consecutive peaks (peak span value is 2). Note how going one more peak on the left side of the mass spectrum increments the protein charge by one unit. Also, the  $M_r$  does not change significantly. Of course, it should be identical in both cases, but that requires zooming-in, enlarging at maximum the spectrum region and carefully positioning the cursor both at the start and end peaks.

FIGURE 4.2: CHARGE STATE-BASED MASS DECONVOLUTION

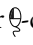
The status bar of the window documents the current inter-peak distance measurement operation that is performed by  $\odot$ -click-drag of the cursor from the start peak to the end peak. The start peak is marked with a green marker and the end peak is marked with a red marker. Start and end positions are documented in the form [left m/z—right m/z] (even if the mouse drag was from right to left, the values are sorted). The m/z delta value documents the distance between both start and end positions. When the end position matches a theoretically expected distance corresponding to a charge difference of 1 or more (depending on the value of the peak span), then the charge  $z$  of the peak under the cursor is provided and the molecular mass ( $M_r$ ) is provided for the analyte whose peak is under the cursor (m/z value documented).



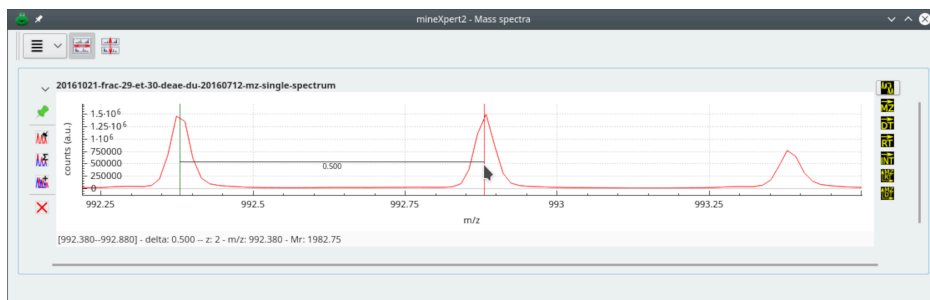
## NOTE

Note that the charge calculation almost never produces an integer value with no fractional part (say, charge  $z=15.0$ ) because it is almost impossible to drag the mouse cursor the exact  $m/z$  range that would lead to such an integral charge value. Almost always, the charge that is calculated looks like 14.995 or 15.001, for example. Why is it impossible to drag the mouse cursor exactly the interval that would produce an integral charge value? Simply because the mouse moves at discrete positions on the screen and these positions might be more or less far apart, depending on the mouse capabilities and on the current zoom factor over the mass spectrum region of interest. It is advised to zoom-in as much as possible over the peaks at hand so as to minimize the difficulties above. It may happen, however, that even zoomed-in peaks are not sufficiently distant to allow a charge calculation. In this case, reduce the stringency over the fractional part that is allowed in the charge (see menu item *Set charge minimal fractional part* at [FIGURE 4.1, “MASS SPECTRUM PLOT WIDGET-SPECIFIC DECONVOLUTION MENU”](#)). By default, the stringency is set at 0.99, that is, any calculated value that has a fractional part either superior or equal to 0.99 or inferior or equal to 0.01 would lead to a successful round-up/round-down to the nearest integer value. Outside of the [0.01-0.99] interval, no charge calculation is performed and thus no deconvolution is performed. When the stringency is too high, reducing it will allow the deconvolution to be carried-over. My own experience is that setting that value to 0.995 is fine for most situations and provides very reliable results.

## 4.2 DECONVOLUTION BASED ON ISOTOPIC CLUSTER PEAKS

In this kind of deconvolution, the user -click-drags the cursor between the first two peaks (when possible) of the isotopic cluster. The charge state of the ion is the inverse of the distance between the two consecutive peaks (that is, the  $m/z$  delta value). [FIGURE 4.3, “ISOTOPIC CLUSTER-BASED MASS DECONVOLUTION”](#) shows that deconvolution process at work.





The user has performed a click-drag movement between the peak under the green marker and the peak under the red marker. The  $m/z$  distance between the two markers is computed and the inverse is the charge of the analyte under this isotopic cluster. If the first and second left peaks of the cluster are not suitable for easy centroid location positioning of the mouse cursor, use peaks to the right and remove the corresponding number of bumps to the right from the calculated monoisotopic mass value.

**FIGURE 4.3: ISOTOPIC CLUSTER-BASED MASS DECONVOLUTION**

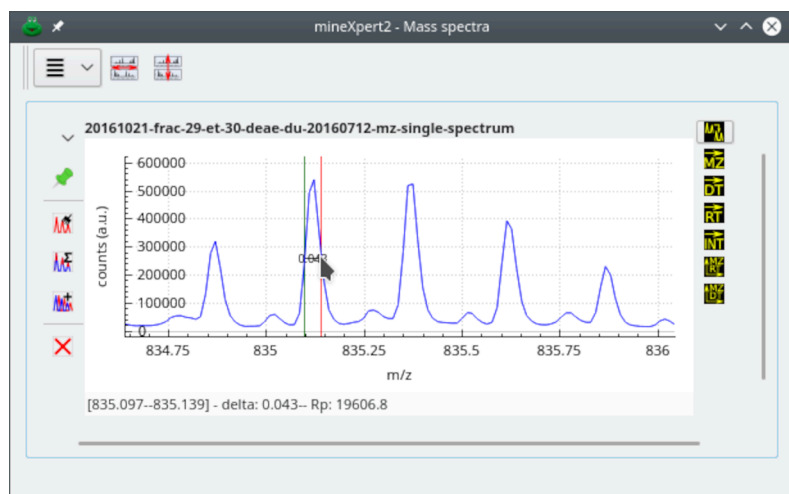


## NOTE: THE MOUSE DRAG POSITION IS NOT SIGNIFICANT

Note that the  $\mathbb{Q}$ -click-dragging direction (left→right or right→left) has no impact on the value of monoisotopic mass computed because the software postulates that the lightest ion is the peak on the left.

## 4.3 READING THE RESOLVING POWER BASED ON MASS SPECTRAL DATA

When  $\mathbb{Q}$ -click-dragging the mouse cursor between two mass spectrum locations of interest, the program computes the apparent resolving power. This process is shown on **FIGURE 4.4, “CALCULATION OF THE RESOLVING POWER”**, where the resolving power is calculated by dragging the mouse cursor from one edge of a peak to the other at half maximum height (this is called *full width at half maximum* [FWHM] resolution).



Click-dragging the mouse cursor will trigger the calculation of the resolving power of the instrument. That value is printed in the status bar.

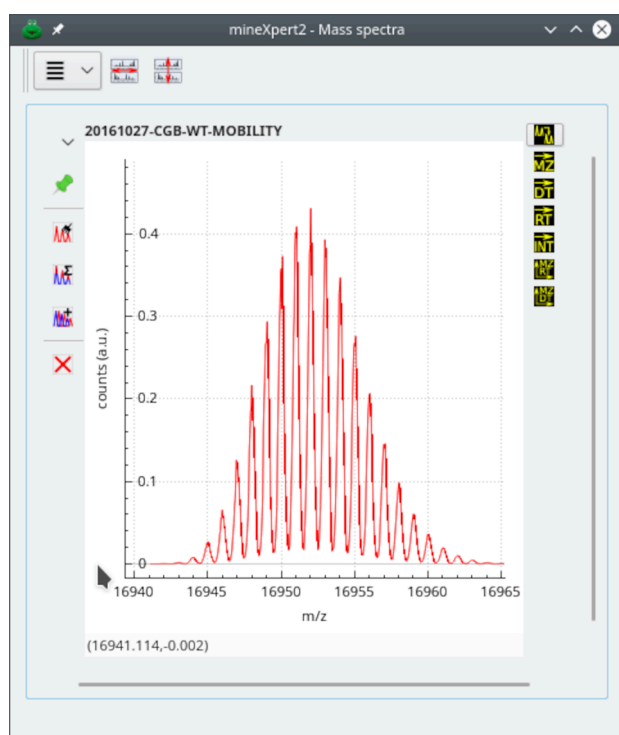
**FIGURE 4.4: CALCULATION OF THE RESOLVING POWER**

## 5 ISOTOPIC CLUSTER CALCULATIONS

### 5.1 CALCULATING ISOTOPIC CLUSTERS WITH ISOPEC

When the resolution of the mass spectrometer is good, zooming-in on a mass peak may reveal that a given ion has given rise not to one peak but to a set of peaks. This set of peaks is called an “*isotopic cluster*”.

It is possible to predict how a given ion (of known chemical formula) is supposed to be revealed in a mass spectrum, in the form of such an isotopic cluster. One such cluster is shown in **FIGURE 5.1, “CALCULATION OF THE ISOTOPIC CLUSTER OF AN ANALYTE”**, for the horse apomyoglobin protein in its monoprotonated form  $[M+H]^+$ , of elemental composition  $C_{769}H_{1213}N_{210}O_{218}S_2$  (this formula is typeset like this intentionally, to show how the formulæ need to be entered in the IsoSpec module).



Calculated isotopic cluster for the monoprotonated apomyoglobin protein.

**FIGURE 5.1: CALCULATION OF THE ISOTOPIC CLUSTER OF AN ANALYTE**

mineXpert2 provides an interface to the `libIsoSpec++` library.

IsoSpec: Hyperfast Fine Structure Calculator

Mateusz K. Łącki, Michał Startek, Dirk Valkenborg, and Anna Gambin

*Analytical Chemistry*, 2017, 89, 3272–3277

DOI: 10.1021/acs.analchem.6b01459

This library performs high-resolution isotopic cluster calculations. In order to run the calculations, it is necessary to have the following items ready:

- An elemental composition formula of the analyte (for example,  $\text{H}_2\text{O}_1$ ). This formula needs to account for the ionization agent that is involved in the ionization of the analyte prior to its detection in the mass spectrometer.



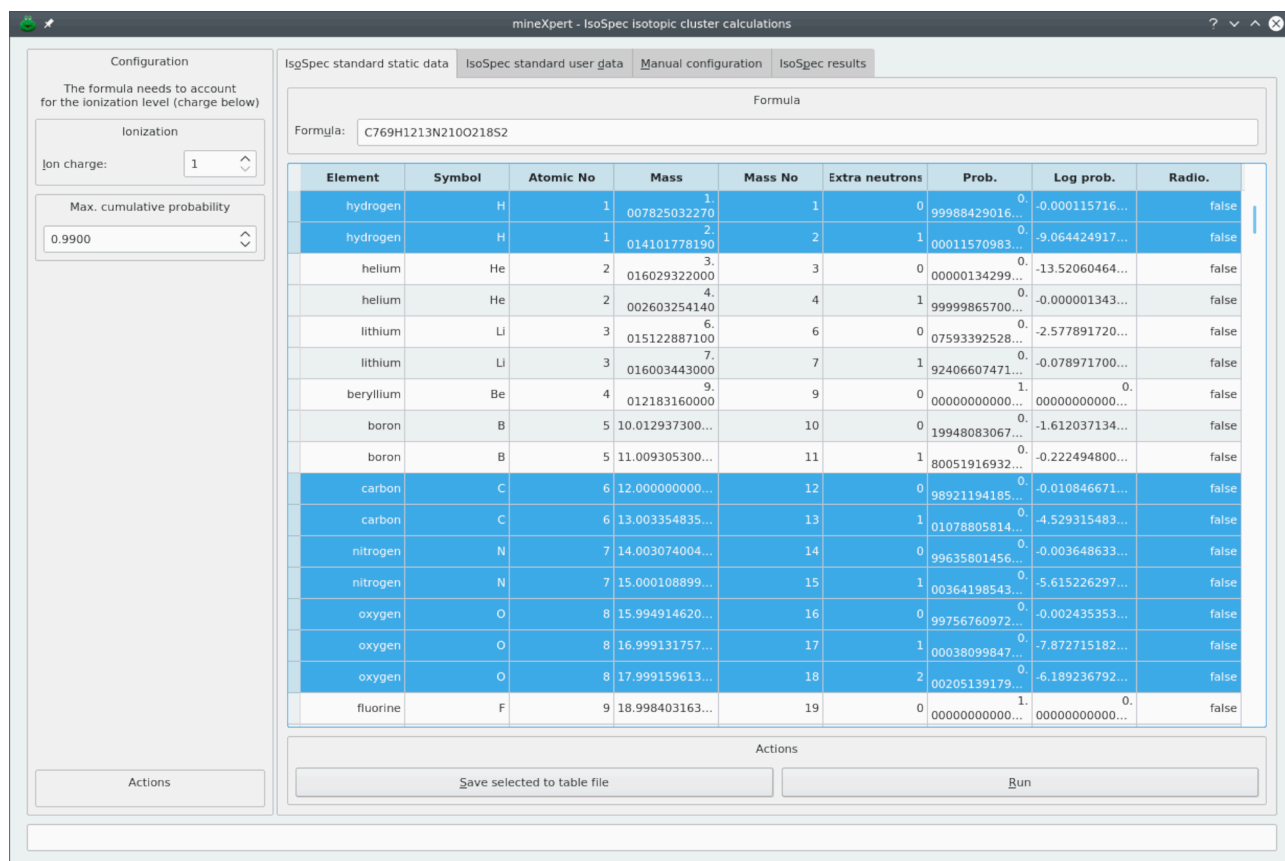
### IMPORTANT

The IsoSpec software requires that all the chemical elements of a chemical formula be indexed. This means that, for water, for example, the formula should be  $\text{H}_2\text{O}_1$  (notice the index 1 after the O element symbol).

- A detailed isotopic configuration of all the chemical elements that are used in the elemental composition formula. mineXpert2 provides two interfaces to define the isotopic characteristics of the chemical elements. These will be described in the following sections.

#### 5.1.1 THE ISOPEC GRAPHICAL USER INTERFACE IN MINEXPRT

Generating isotopic clusters using the IsoSpec software package is not easily carried over, in particular because this remarkable library is designed to be highly performant. The authors rightfully put their energy into optimizations for accuracy and speed instead of investing in a graphical user interface. mineXpert2 provides that graphical user interface, shown in **FIGURE 5.2**, “ISOTOPIC CLUSTER CALCULATION DIALOG WINDOW”, that shows up upon selection of the program's main window's *Utilities > Isotopic cluster calculations* menu.



The dialog window contains two panels. The left hand side panel configures the charge for which the calculation is to be carried over and the maximum cumulative isotopic presence probability that IsoSpec must reach during the calculation. The right hand side panel contains a tab widget that contains the configuration tabs and the results tab.

**FIGURE 5.2: ISOTOPIC CLUSTER CALCULATION DIALOG WINDOW**

An isotopic cluster calculation is most probably performed with the aim of simulating an expected isotopic cluster for an analyte that is being analyzed by mass spectrometry. It is thus logical that the analyte be in an ionized form. The way that the analyte has been ionized needs to be taken into account in the chemical formula that describes the ion for which the isotopic cluster is being calculated. For example, when determining the chemical formula of a protein in the positive ion mode, the number of protons used to ionize the protein need to be included in the analyte elemental composition formula.



## WARNING

The IsoSpec software is “charge-agnostic” in the sense that it does not know what element in the chemical formula is responsible for the ionization of the analyte. Therefore, IsoSpec does not know of (and does not care about) the charge of the analyte. The ionization level of the analyte can be handled by mineXpert2 if that information is set to the *Ion charge* spin box widget. By default, the charge state of the analyte is 1.

The *Max. cumulative probability* spin box widget serves to configure the extent to which IsoSpec simulates the theoretically expected isotopic cluster. A value of 0.99 tells the software to simulate enough combinations of the analyte isotopes to represent 99 % of the theoretically expected combinations.



## TIP

For large biopolymers, it might be prudent to start with a relatively low value for *Max. cumulative probability*, because setting this value too high, that is, near 1, would increase notably the calculation duration.

To perform isotopic cluster calculations, the simulation software needs to be aware of all the isotopes of all the chemical elements that enter in the composition of the ionized analyte. An isotope is defined by its mass and by the probability that it is found in nature. Carbon has two major isotopes that can be found in nature: the [<sup>12</sup>C] most abundant isotope and the [<sup>13</sup>C] least abundant isotope (the [<sup>14</sup>C] isotope is irrelevant for conventional mass spectrometry unless molecules have been artificially enriched in that isotope).

There are three ways that the user might use the IsoSpec software package. Two of them involve a configuration preparation on the part of the user. The third one, that we'll describe first, does not involve any chemical element configuration. The other ways are reviewed in the next sections.

### 5.1.1.1 STATIC STANDARD ELEMENT TABLES SHIPPED WITHIN THE ISOPEC LIBRARY

In order to document all the chemical elements' isotopes' characteristics, the IsoSpec library has, in its own headers, a number of arrays that mineXpert2 automatically loads up when opening the **FIGURE 5.2, “ISOTOPIC CLUSTER CALCULATION DIALOG WINDOW”** dialog window.<sup>1</sup> These data are displayed in the *IsoSpec standard static data* tab. The table view widget is not editable, hence the *static* qualifier in the tab name.

In this tab, all the user has to do is enter the formula for which the isotopic cluster calculation is to be performed. The formula needs to be pasted in the *Formula* line edit widget.

---

<sup>1</sup> That process is performed at build time and is thus not configurable at run time.

Once the formula has been set to its line edit widget, configure the charge of the ion (*Ion charge* spin box widget) and the *Max. cumulative probability*. Now, click *Run*.

#### 5.1.1.2 MODIFIED STANDARD ELEMENT TABLES SHIPPED WITHIN THE ISOPEC LIBRARY

While the previous section showed how to use the static element tables from the IsoSpec library, this section shows a way to configure these elemental data. For this, start from the standard IsoSpec-based data and modify them according to one's needs. To proceed, select all the element rows of interest from the *IsoSpec standard static data* tab and click *Save selected to table file*.



#### TIP

The order with which the rows are selected is respected in the export, so make sure to select the rows in the order that makes sense for you. Also, be sure to select rows and not only some cells; to achieve, this click in the left margin of the table view widget.

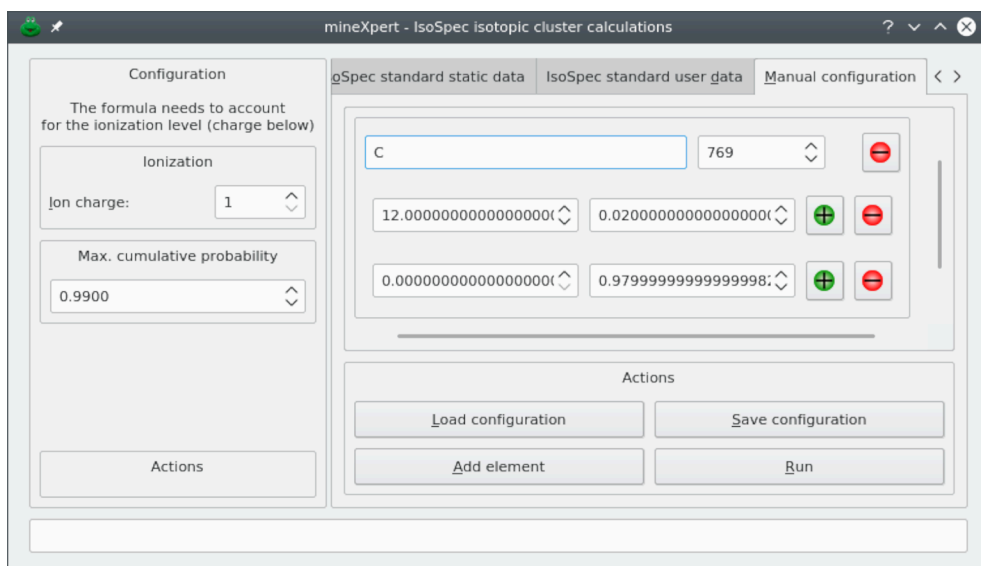
The export is performed as a text TSV (tab-separated value format) file in a layout that closely mimicks the data visible in the table view. Once the file has been saved, open it up in LibreOffice and modify the data to suit your needs. For example, in case of a labelling with [<sup>13</sup>C], with an efficiency of almost 98 %, let's change the [<sup>12</sup>C] abundance (probability) to 0.02 and the [<sup>13</sup>C] abundance to 0.98.

Now that the modifications were performed, save the file under the same format and load it in the *IsoSpec standard user data* tab of the dialog window by clicking *Load table from file*. The table view now shows the new carbon abundance values. From now on, use this tab exactly like you did for the standard isospec *static* data work described above.

Here also, it is possible to select determinate rows and then save them to a file for later use.

#### 5.1.1.3 USER HAND-MADE CONFIGURATION OF THE ELEMENT DATA WITHIN MINEXPRT

The last way the user might configure the chemical elements to be used during an isotopic cluster calculation is based on the fully manual description of the elements and of the isotopes. That configuration is performed in the *Manual configuration* tab of the dialog window. This method is slightly more involved than the previous one but provides also for a much greater flexibility: it allows one to create “new chemical elements” that might be required in specific labelling experiments. The manual configuration is carried over in the *Manual configuration* tab of the dialog, as shown in **FIGURE 5.3, “HAND-MADE USER CONFIGURATION OF THE CHEMICAL ELEMENTS AND FORMULA”**.



When the dialog is created, the tab is empty. To start creating element definitions, click *Add element*.

**FIGURE 5.3: HAND-MADE USER CONFIGURATION OF THE CHEMICAL ELEMENTS AND FORMULA**

Upon creation of the dialog window, the *Manual configuration* tab is empty, with only two rows of buttons at the bottom of the tab. To start configuring chemical elements, click *Add element* to create an “element group box” that contains a number of widgets organized in two rows:

- Top row, a line edit widget to receive the chemical element symbol, *C* in the example;
- A spin box widget in which to set the number of such atoms in the formula for which the isotopic cluster is being calculated. In the example, we set this value to *769*;
- A button with a “minus” image that removes all the “element group box” in one go;
- The bottom row contains an “isotope frame widget” with two spin boxes for the mass of the isotope being configured (left) and its corresponding abundance (right);
- In addition to the spin boxes, two buttons, with a “plus” or a “minus” figure, allow one to respectively add or remove isotope frames.



## NOTE

It is not possible to remove all the isotope frames from an element group box, otherwise that group box would become useless.

Once an isotope frame has been filled-up, a new line might be required. To create a new isotope frame widget, click any “plus”-labelled button in any of the isotope frames. Once a new frame is created, the spin box widgets that it contains are set to *0.00000*. Fill-in these spin boxes with mass and abundance and go on along this path to create as many isotopes as required.



Once all the isotopes for a given chemical element have been defined, a new element might be needed. For this, click *Add element* and start the configuration of the new element as described above.

The manual isotopic configuration of the chemical elements required to perform an isotopic cluster calculation for a given formula is tedious. The user may want to save a given configuration to a file (click *Save configuration*) so that it is easier to recreate automatically all the widgets upon loading of that saved configuration (click *Load configuration*).

The final configuration is shown in **FIGURE 5.4, “TYPICAL MANUAL CONFIGURATION OF THE ISOTOPIC CHARACTERISTICS OF THE CHEMICAL ELEMENTS”**. The experiment that was configured above is a labelling of a glucose molecule with Cz, an imaginary chemical element that is like carbon but that has a [ $^{14}\text{C}$ ]. The glucose molecule (normal formula:  $\text{C}_6\text{H}_{12}\text{O}_6$ ) is labelled on one single carbon atom with an efficiency of 95 %. This means that, when the labelling fails (in 5 % of the cases) the carbon atom has its isotopes with usual probabilities (compounded by the fact that the normal atom is found at that position only in 5 % of the cases). The isotopic abundances for the Cz element are thus:

- For [ $^{12}\text{C}$ ]:  $0.05 * \text{normal } [^{12}\text{C}] \text{ abundance}$ ;
- For [ $^{13}\text{C}$ ]:  $0.05 * \text{normal } [^{13}\text{C}] \text{ abundance}$ ;
- For [ $^{14}\text{C}$ ]: 0.95;

The screenshot shows the 'mineXpert - IsoSpec isotopic cluster calculations' window. The 'Manual configuration' tab is selected. On the left, the 'Configuration' panel has a note: 'The formula needs to account for the ionization level (charge below)'. It shows 'Ionization' with 'Ion charge: 1' and 'Max. cumulative probability' set to '0.9900'. The main area lists elements: C (count 5), H (count 12), O (count 6), and Cz (count 1). Each element has a table of isotopes with their natural abundances and user-defined relative abundances. For example, Carbon (C) has two isotopes: 12.000000000000000 (natural 0.98921194185046690, user 0.95) and 13.0033548351999996 (natural 0.01078805814953308, user 0.05). Hydrogen (H) has two isotopes: 1.00782503206999995 (natural 0.99988429016430790, user 0.99) and 2.01410177780000010 (natural 0.00011570983569203, user 0.01). Oxygen (O) has three isotopes: 15.9949146195599993 (natural 0.99756760972956104, user 0.99), 16.9991316999999995 (natural 0.00038099847600609, user 0.01), and 17.9991610000000008 (natural 0.00205139179443282, user 0.00). Carbon-13 (Cz) has three isotopes: 12.000000000000000 (natural 0.04946060000000000, user 0.05), 13.0033548351999996 (natural 0.00053940290747665, user 0.95), and 14.0032419889999992 (natural 0.9499999999999995, user 0.00). At the bottom, the 'Actions' panel contains buttons for 'Load configuration', 'Save configuration', 'Add element', and 'Run'.

The user has configured a labelling experiment where the glucose molecule is labelled at a single carbon position with a [ $^{14}\text{C}$ ] atom (the efficiency of the labelling is 95 %).

**FIGURE 5.4: TYPICAL MANUAL CONFIGURATION OF THE ISOTOPIC CHARACTERISTICS OF THE CHEMICAL ELEMENTS**



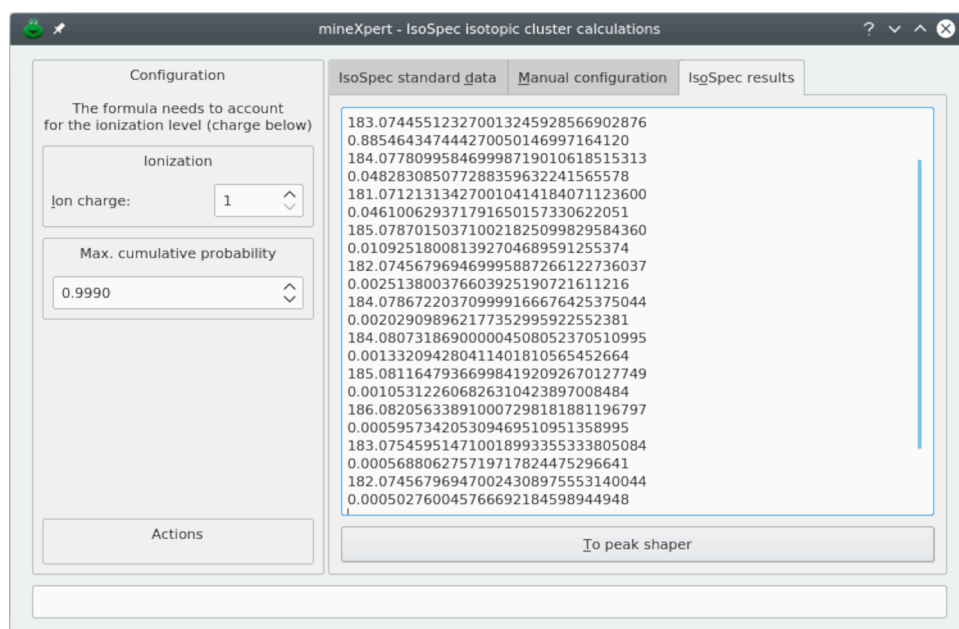
## NOTE

Note that the “normal” carbon count is 5 (and not 6), that the hydrogen count is 13 (and not 12, because the glucose is protonated) and the labelled carbon is present only once.

#### 5.1.1.4 THE ISOspec RESULTS ARE NOT SHAPED MASS PEAKS

Once the configurations have been terminated, the calculations can finally be performed by the IsoSpec library. In the manual configuration setting, the formula is automatically handled since each chemical element that is defined goes along with the count of the corresponding atoms. In the case of the standard IsoSpec configuration (either modified or not), the user has to enter the chemical formula of the analyte in the *Formula* line edit widget. Click *Run*. If the configuration was correct and IsoSpec could run the calculation properly, then the dialog window switches to the *IsoSpec results* tab (FIGURE 5.5, “RESULTS FROM THE ISOTOPIC CLUSTER CALCULATION”). That tab contains a text edit widget in which the results are displayed.

Note that the  $m/z$  values calculated by IsoSpec are “corrected” for the charge level that was specified in the left panel of the dialog window prior to their display in the results tab (FIGURE 5.2, “ISOTOPIC CLUSTER CALCULATION DIALOG WINDOW”).



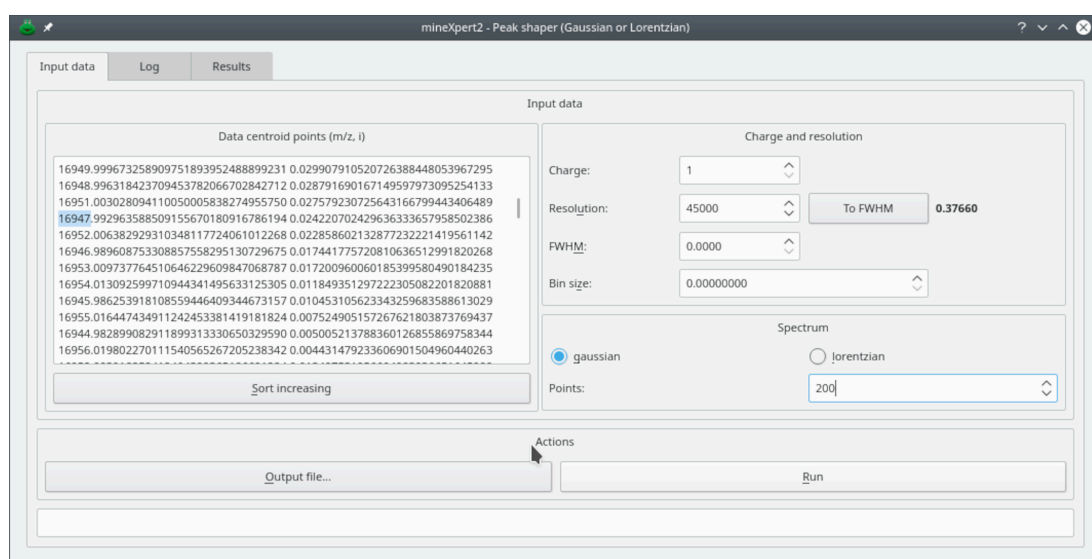
The IsoSpec library computes the probability of the various combinations of all the isotopes that make the chemical formula submitted to it. The results are in the form of peak centroid values along with corresponding probabilities. The sum of the probabilities corresponds to the *Max. cumulative probability* value that was set by the user.

FIGURE 5.5: RESULTS FROM THE ISOTOPIC CLUSTER CALCULATION

The results that are produced by IsoSpec represent the peak centroids of the isotopic cluster. The results are thus a set of ( $m/z$ ,  $i$ ) pairs that have not the characteristic shape (the profile) that is found in mass spectra. mineXpert2 features the ability to give a shape to the centroid peaks. For that, click *To peak shaper* to open the *Peak Shaper* dialog window preloaded with the IsoSpec-generated peak centroids. The workings of this peak shaping feature is described in SECTION 5.2, “SHAPING MASS PEAK CENTROIDS INTO WELL-PROFIED PEAKS”.

## 5.2 SHAPING MASS PEAK CENTROIDS INTO WELL-PROFILED PEAKS

The shape of mass peaks is typically Gaussian or Lorentzian (or a mix thereof). There are some data simulation or analysis processes that lead to having mass peaks characterized by a single centroid  $m/z$  value and a corresponding intensity. Plotted to a graph, a centroid mass peak yields a bar. In order to convert mass peak centroids into something that resembles a real “profile” mass peaks, a mathematical formula can be applied, with some parameters to configure the shapes generated. mineXpert2 includes that feature, accessible *via* the *Utilities > Mass peak shape calculations* menu item. The window that opens up is shown in **FIGURE 5.6, “SETTING-UP OF THE CENTROID MASS PEAK SHAPING PROCESS”**



This dialog window allows one to configure the shaping of mass centroid peaks. Setting a spectrum name in the *Mass spectrum name* line edit widget will help recognize the result mass spectrum once displayed in the *Mass spectra* window (see below).

**FIGURE 5.6: SETTING-UP OF THE CENTROID MASS PEAK SHAPING PROCESS**

The mass centroid peaks are listed in the *Data centroid points ( $m/z, i$ )* text edit widget. These values are pasted there by the user or copied automatically from the isotopic cluster calculation dialog window (see **SECTION 5.1.1.4, “THE ISO SPEC RESULTS ARE NOT SHAPED MASS PEAKS”**). The width of the “profile” mass peak is determined either by setting the resolution of the instrument (in the example, that is set to 45000) or by setting the width of the peak at half maximum of its height (FWHM).

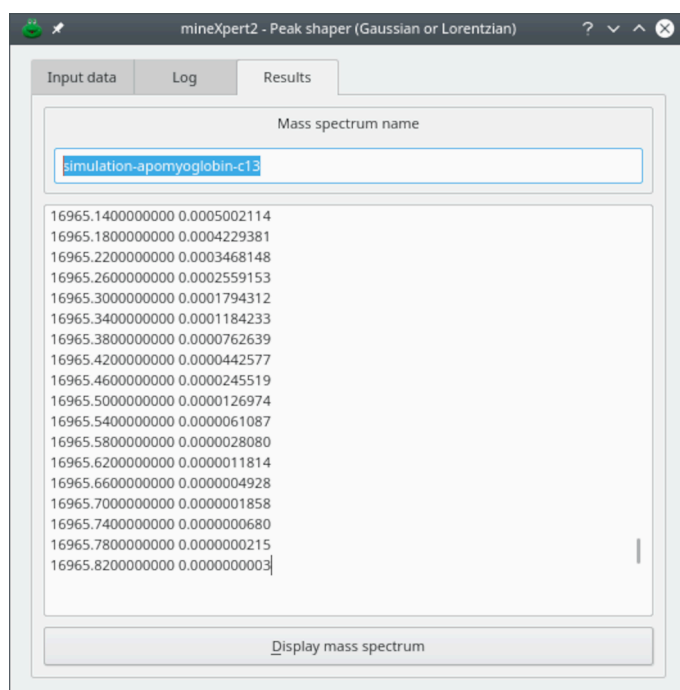
The spectrum that is generated can be of a Gaussian or a Lorentzian shape. That parameter is configured by selecting the corresponding radio button widget. The number of points used to actually craft the shape of the peak is configurable. In the example, that parameter is set to 200.



## TIP

When defining the parameters for the peak shaping process, it might be useful to have an idea of the FWHM value at a given  $m/z$  value for a given resolution. To compute that FWHM value, double-click-select a single mass peak centroid  $m/z$  value from the text edit widget, set the resolving power of your instrument and then click *To FWHM*. The computed FWHM value will be displayed next to the button. The text is selectable so as to copy it to the clipboard and then in the *FWHM* spin box widget.

Once the peak shaping parameters have been set, click *Run*. The dialog window shifts tab to *Results*, as shown in **FIGURE 5.7, “SHAPED PEAKS MASS SPECTRAL DATA”**.



The mass spectral data corresponding to the combination of each individual “peak-shaped” centroid peak are displayed in this tab. Plotting these data will produce a nicely shaped isotopic cluster matching what would be seen by recording a mass spectrum on an instrument.

**FIGURE 5.7: SHAPED PEAKS MASS SPECTRAL DATA**



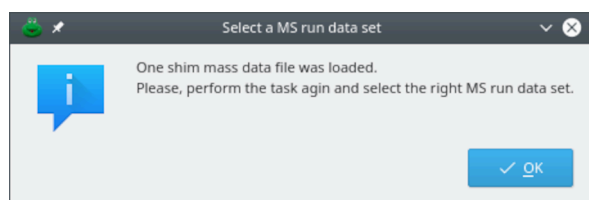
## TIP

To ensure that various isotopic cluster calculation-based mass spectra can later be recognized, the user is advised, for the different simulations, to insert distinct names in the *Mass spectrum name* line edit widget. This name will be used later when creating the mass spectrum if the user asks that the mass spectrum

be displayed by clicking *Display mass spectrum*. In the eventuality that the mass spectrum name is not changed from a simulation to the other, a safeguard process ensures that names are absolutely unambiguous by appending to the mass spectrum name the time at which the mass spectrum is displayed.

It is possible to automatically plot the result mass spectrum by clicking *Display mass spectrum*. The process gets a little involved here because the mass spectrum is not created from a file, and mineXpert2's handling of plots is based on the fact that a MS run data set has been created, typically at data file loading time. Two situations may arise:

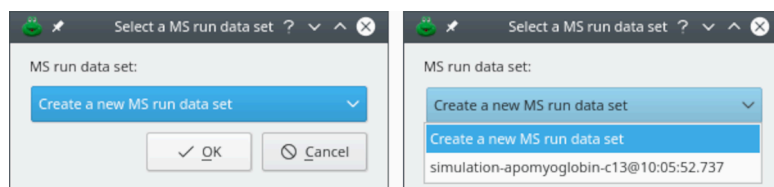
- If not a single MS run data set is currently available (that is, no single mass spectrometry data file was loaded during the current session), the program will inform the user that a shim MS run data set was created (FIGURE 5.8, “CREATION OF A SHIM MS RUN DATA SET”). At this time, the dialog window listing all the open MS run data sets will show up (FIGURE 2.2, “THE OPEN MS RUN DATA SETS”). The information message advises to perform the task again, that is, click *Display mass spectrum*.



This information dialog window informs the user that a new shim MS run data set was created. Shim MS run data sets need to be created when new plots are created for data that were not loaded from a disk file. The user is advised to repeat the last action. In this case, click *Display mass spectrum* again. The next step is described in the text below.

**FIGURE 5.8: CREATION OF A SHIM MS RUN DATA SET**

- If at least one MS run data set was loaded already in the program, then an input dialog window shows up letting the user select the MS run data set to which the isotopic cluster calculation mass spectrum is to be anchored to (FIGURE 5.9, “SELECTION OF A MS RUN DATA SET”). It is always possible to ask that a new MS run data set be created *ex nihilo* for the new mass spectrum by selecting *Create a new MS run data set*. In this case, the process goes back to the previous item. This process is useful if each mass spectrum must have a distinct color.



If one or more MS run data set(s) preexisted in the program, this input dialog shows up (left). Upon clicking the combo box widget, the user can select the MS run data set to anchor the new mass spectrum to (right).

**FIGURE 5.9: SELECTION OF A MS RUN DATA SET**

The isotopic cluster calculation-based mass spectrum shows up in the *Mass spectra* window, as shown in **FIGURE 5.10, “SPECTRUM CREATED USING THE PEAK SHAPING FEATURE”**.



The spectrum corresponds to a combination of each individual spectrum obtained by shaping each individual mass peak centroid in the input data list.

**FIGURE 5.10: SPECTRUM CREATED USING THE PEAK SHAPING FEATURE**

Note that if the resolution asked is very high, the resulting shaped mass peaks might appear a bit “hairy”. By tweaking the *Bin size* value, the binning of the spectra might improve the situation. Otherwise, using the plot widget main menu to apply a Savitzky-Golay filter, described at **SECTION 3.1.6, “SAVITZKY-GOLAY FILTERING OF ANY KIND OF DATA” (PAGE 40)**, will certainly improve things.



## TIP

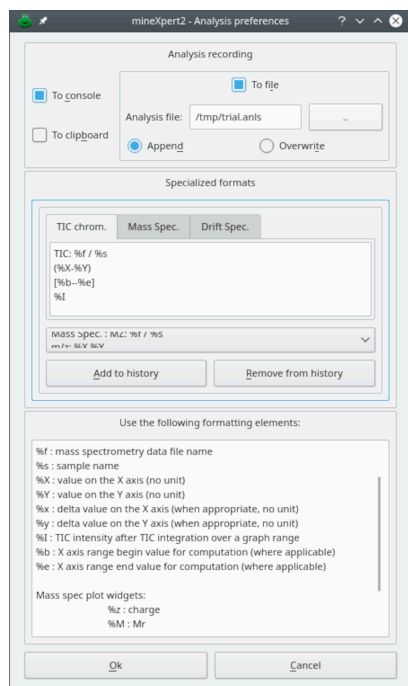
If the peak centroids were not “corrected” for their charge in the previous generation step (as in the case of the isotopic cluster calculation, it is still time to apply this “correction” by setting the charge in the *Charge* spin box widget. If the charge was already accounted for, as described in [SECTION 5.1.1.4, “THE ISOPEC RESULTS ARE NOT SHAPED MASS PEAKS” \(PAGE 60\)](#), then leave the charge to 1 and the results will be correct.



## 6 RECORDING DATA MINING DISCOVERIES

When doing mass analysis work it is often desirable to store the painstakingly manually picked  $m/z$  or  $M_r$  values for later use. mineXpert2 provides a flexible process allowing one to record the data mining discoveries to a file, the clipboard, the console or any combination thereof.

In order to record the innumerable analysis steps that make a data mining session, the *File > Analysis preferences...* menu will open a dialog window shown in **FIGURE 6.1**, “SETTING-UP OF THE RECORDING OF THE DATA MINING DISCOVERIES”.



It is possible to configure the recording system to record to either the console, the clipboard, a file (in append mode or in overwrite mode) or any combination thereof. The format of the string is defined using special characters (see text) and might be defined specifically for the three main graphs: TIC/XIC chromatogram, mass spectrum and drift spectrum.

**FIGURE 6.1: SETTING-UP OF THE RECORDING OF THE DATA MINING DISCOVERIES**

In that window, the user can select the destination of the data analysis recording system: console, clipboard, file or any combination of the three. When selecting file recording, the user might specify if the recording would overwrite any preexisting file or, instead, append to that file. Depending on the kind of plot where data mining occurs, the format of the data to be recorded needs to change. Indeed, it would make no sense to record the charge  $z$  when mining data in the *Drift spectra* window. This is why the text format of the data export needs to be defined for each one of the three kinds of plots: TIC/XIC chromatogram, mass spectrum or drift spectrum (*Specialized formats* group box).

Before delving into the configuration intricacies, let us tell immediately how to trigger the recording of the mining discoveries: using the `SPACE` bar in the composite plot widget containing the graph being analyzed.

The format used to define the text string to be stored on console and/or in file can contain particular tokens as described below:

- `%f`: mass spectrometry data file name.
- `%s`: sample name.
- `%X`: value on the x-axis of the graph (no unit). For a drift spectrum, that would be drift times in milliseconds, for a mass spectrum, that would be  $m/z$  values, for a TIC/XIC chromatogram, that would be retention times in minutes.
- `%Y`: value on the y-axis of the graph (no unit). In all the graph plots, that would be intensities in any unit provided by the mass spectrometer (typically, counts).
- `%x`: delta value on the x-axis (when appropriate, no unit).
- `%y`: delta value on the y-axis (when appropriate, no unit).
- `%I`: TIC intensity after TIC integration over a graph range
- `%b`: x-axis range *begin* point for the computation (where applicable, for example for the TIC integration to a single value).
- `%e`: x-axis range *end* point for the computation (where applicable).
- *For mass spec plot widgets:*
  - `%z`: charge
  - `%M`:  $M_r$  as computed during deconvolution (see [CHAPTER 4, MASS SPECTRAL DECONVOLUTIONS](#)).



## NOTE

It is important to keep in mind that the `%z` and `%M` format strings can only work if the user is actually analyzing a mass spectrum and if the user has *effectively* performed a deconvolution operation that has allowed computing these two values. If the values are not available, the program shows *nan* (“not a number”) in the textual output upon hitting the `SPACE` bar (see below). Likewise, the `%I` format string

will only hold meaningful data if an integration to a single TIC intensity value has been determined (see [FIGURE 3.10, “INTEGRATING MASS DATA FROM A MASS SPECTRUM TO A SINGLE TIC INTENSITY VALUE”](#)) for that specific data analysis record.

The recording of data analysis steps works in any trace plot widget (not yet in the color map plot widgets) even if there are more than one trace in a given plot widget. In the case the plot widget has more than one trace and none or more than one is selected, the program will ask that a single trace be selected. If data are to be scrutinized for more traces, simply select each trace in turn and trigger the analysis stanza to be output using the `SPACE` bar each time.

Once configured, the format strings might be stored in a drop down box for later use. To that end, click onto the *Add to history* button while having the format text displayed in the text editor and it will be appended to the drop-down list. The list gets stored when the dialog window is closed and will be filled-up again when the program is restarted.

As an example, if the user defined the following format string for a mass spectrum graph:

```
Mass spec. :  
mz = (%X, %Y) z = %z  
filename = %f  
date = 20161021  
session = 20161021  
mslevel = 1 msion = esi msanal = tof  
chrom = DEAE fraction = 25  
seq = pos = oxlevel = 0 pos =  
intensity =  
comment =
```

Then, a resulting data mining stanza that would be recorded will look like this:

```
Mass Spec. :  
mz = (1051.8, 50863) z = 1  
filename = 20161017-rusconi-frac-25-deae-20160712.mzml  
date = 20161021  
session = 20161021  
mslevel = 1 msion = esi msanal = tof
```

```
chrom = DEAE fraction = 25  
seq = pos =   oxlevel = 0 pos =  
intensity =  
comment =
```

Interestingly, the user can define any kind of format, leaving fields available for later filling-in. This feature is of immense value when the analysis file is used later to fill-in a database for easy storage and interrogation of the mining discoveries. In this case, it would be useful to have the file opened in an editor and at each new stanza, edit the *comment* field if something needs to be commented, like the shape/intensity of a mass peak, for example.

Note that the program closes the file each time a new stanza has been written. This makes it possible to edit that file safely in between each stanza record. Remember to force the editor to reload the file from disk after each mining discovery recording.

When the recording involves sending the analysis data to the console, the data are sent to it as text that is colored the same as the spectrum that was under scrutiny.

When the mouse cursor has been placed at the proper location on the graph (with or without  $\mathbb{Q}$ -click-dragging, depending on the situation), the user hits the SPACE bar and the data analysis stanza is recorded to the selected destination(s): console, clipboard, file.

# A GNU GENERAL PUBLIC LICENSE VERSION 3

Version 3, 29 June 2007

Copyright © 2007 Free Software Foundation, Inc. [HTTPS://FSF.ORG/](https://fsf.org/) 

Everyone is permitted to copy and distribute verbatim copies of this license document, but changing it is not allowed.

## PREAMBLE

The GNU General Public License is a free, copyleft license for software and other kinds of works.

The licenses for most software and other practical works are designed to take away your freedom to share and change the works. By contrast, the GNU General Public License is intended to guarantee your freedom to share and change all versions of a program—to make sure it remains free software for all its users. We, the Free Software Foundation, use the GNU General Public License for most of our software; it applies also to any other work released this way by its authors. You can apply it to your programs, too.

When we speak of free software, we are referring to freedom, not price. Our General Public Licenses are designed to make sure that you have the freedom to distribute copies of free software (and charge for them if you wish), that you receive source code or can get it if you want it, that you can change the software or use pieces of it in new free programs, and that you know you can do these things.

To protect your rights, we need to prevent others from denying you these rights or asking you to surrender the rights. Therefore, you have certain responsibilities if you distribute copies of the software, or if you modify it: responsibilities to respect the freedom of others.

For example, if you distribute copies of such a program, whether gratis or for a fee, you must pass on to the recipients the same freedoms that you received. You must make sure that they, too, receive or can get the source code. And you must show them these terms so they know their rights.

Developers that use the GNU GPL protect your rights with two steps: (1) assert copyright on the software, and (2) offer you this License giving you legal permission to copy, distribute and/or modify it.

For the developers' and authors' protection, the GPL clearly explains that there is no warranty for this free software. For both users' and authors' sake, the GPL requires that modified versions be marked as changed, so that their problems will not be attributed erroneously to authors of previous versions.

Some devices are designed to deny users access to install or run modified versions of the software inside them, although the manufacturer can do so. This is fundamentally incompatible with the aim of protecting users' freedom to change the software. The systematic pattern of such abuse occurs in the area of products for individuals

to use, which is precisely where it is most unacceptable. Therefore, we have designed this version of the GPL to prohibit the practice for those products. If such problems arise substantially in other domains, we stand ready to extend this provision to those domains in future versions of the GPL, as needed to protect the freedom of users. Finally, every program is threatened constantly by software patents. States should not allow patents to restrict development and use of software on general-purpose computers, but in those that do, we wish to avoid the special danger that patents applied to a free program could make it effectively proprietary. To prevent this, the GPL assures that patents cannot be used to render the program non-free.

The precise terms and conditions for copying, distribution and modification follow.

## TERMS AND CONDITIONS

### 0. DEFINITIONS.

“This License” refers to version 3 of the GNU General Public License.

“Copyright” also means copyright-like laws that apply to other kinds of works, such as semiconductor masks.

“The Program” refers to any copyrightable work licensed under this License. Each licensee is addressed as “you”.

“Licensees” and “recipients” may be individuals or organizations.

To “modify” a work means to copy from or adapt all or part of the work in a fashion requiring copyright permission, other than the making of an exact copy. The resulting work is called a “modified version” of the earlier work or a work “based on” the earlier work.

A “covered work” means either the unmodified Program or a work based on the Program.

To “propagate” a work means to do anything with it that, without permission, would make you directly or secondarily liable for infringement under applicable copyright law, except executing it on a computer or modifying a private copy. Propagation includes copying, distribution (with or without modification), making available to the public, and in some countries other activities as well.

To “convey” a work means any kind of propagation that enables other parties to make or receive copies. Mere interaction with a user through a computer network, with no transfer of a copy, is not conveying.

An interactive user interface displays “Appropriate Legal Notices” to the extent that it includes a convenient and prominently visible feature that (1) displays an appropriate copyright notice, and (2) tells the user that there is no warranty for the work (except to the extent that warranties are provided), that licensees may convey the work under this License, and how to view a copy of this License. If the interface presents a list of user commands or options, such as a menu, a prominent item in the list meets this criterion.

## I. SOURCE CODE.

The “source code” for a work means the preferred form of the work for making modifications to it. “Object code” means any non-source form of a work.

A “Standard Interface” means an interface that either is an official standard defined by a recognized standards body, or, in the case of interfaces specified for a particular programming language, one that is widely used among developers working in that language.

The “System Libraries” of an executable work include anything, other than the work as a whole, that (a) is included in the normal form of packaging a Major Component, but which is not part of that Major Component, and (b) serves only to enable use of the work with that Major Component, or to implement a Standard Interface for which an implementation is available to the public in source code form. A “Major Component”, in this context, means a major essential component (kernel, window system, and so on) of the specific operating system (if any) on which the executable work runs, or a compiler used to produce the work, or an object code interpreter used to run it.

The “Corresponding Source” for a work in object code form means all the source code needed to generate, install, and (for an executable work) run the object code and to modify the work, including scripts to control those activities. However, it does not include the work's System Libraries, or general-purpose tools or generally available free programs which are used unmodified in performing those activities but which are not part of the work. For example, Corresponding Source includes interface definition files associated with source files for the work, and the source code for shared libraries and dynamically linked subprograms that the work is specifically designed to require, such as by intimate data communication or control flow between those subprograms and other parts of the work.

The Corresponding Source need not include anything that users can regenerate automatically from other parts of the Corresponding Source.

The Corresponding Source for a work in source code form is that same work.

## 2. BASIC PERMISSIONS.

All rights granted under this License are granted for the term of copyright on the Program, and are irrevocable provided the stated conditions are met. This License explicitly affirms your unlimited permission to run the unmodified Program. The output from running a covered work is covered by this License only if the output, given its content, constitutes a covered work. This License acknowledges your rights of fair use or other equivalent, as provided by copyright law.

You may make, run and propagate covered works that you do not convey, without conditions so long as your license otherwise remains in force. You may convey covered works to others for the sole purpose of having them make modifications exclusively for you, or provide you with facilities for running those works, provided that you

comply with the terms of this License in conveying all material for which you do not control copyright. Those thus making or running the covered works for you must do so exclusively on your behalf, under your direction and control, on terms that prohibit them from making any copies of your copyrighted material outside their relationship with you.

Conveying under any other circumstances is permitted solely under the conditions stated below. Sublicensing is not allowed; section 10 makes it unnecessary.

### 3. PROTECTING USERS' LEGAL RIGHTS FROM ANTI-CIRCUMVENTION LAW.

No covered work shall be deemed part of an effective technological measure under any applicable law fulfilling obligations under article 11 of the WIPO copyright treaty adopted on 20 December 1996, or similar laws prohibiting or restricting circumvention of such measures.

When you convey a covered work, you waive any legal power to forbid circumvention of technological measures to the extent such circumvention is effected by exercising rights under this License with respect to the covered work, and you disclaim any intention to limit operation or modification of the work as a means of enforcing, against the work's users, your or third parties' legal rights to forbid circumvention of technological measures.

### 4. CONVEYING VERBATIM COPIES.

You may convey verbatim copies of the Program's source code as you receive it, in any medium, provided that you conspicuously and appropriately publish on each copy an appropriate copyright notice; keep intact all notices stating that this License and any non-permissive terms added in accord with section 7 apply to the code; keep intact all notices of the absence of any warranty; and give all recipients a copy of this License along with the Program.

You may charge any price or no price for each copy that you convey, and you may offer support or warranty protection for a fee.



## 5. CONVEYING MODIFIED SOURCE VERSIONS.

You may convey a work based on the Program, or the modifications to produce it from the Program, in the form of source code under the terms of section 4, provided that you also meet all of these conditions:

- a.** The work must carry prominent notices stating that you modified it, and giving a relevant date.
- b.** The work must carry prominent notices stating that it is released under this License and any conditions added under section 7. This requirement modifies the requirement in section 4 to “keep intact all notices”.
- c.** You must license the entire work, as a whole, under this License to anyone who comes into possession of a copy. This License will therefore apply, along with any applicable section 7 additional terms, to the whole of the work, and all its parts, regardless of how they are packaged. This License gives no permission to license the work in any other way, but it does not invalidate such permission if you have separately received it.
- d.** If the work has interactive user interfaces, each must display Appropriate Legal Notices; however, if the Program has interactive interfaces that do not display Appropriate Legal Notices, your work need not make them do so.

A compilation of a covered work with other separate and independent works, which are not by their nature extensions of the covered work, and which are not combined with it such as to form a larger program, in or on a volume of a storage or distribution medium, is called an “aggregate” if the compilation and its resulting copyright are not used to limit the access or legal rights of the compilation’s users beyond what the individual works permit. Inclusion of a covered work in an aggregate does not cause this License to apply to the other parts of the aggregate.

## 6. CONVEYING NON-SOURCE FORMS.

You may convey a covered work in object code form under the terms of sections 4 and 5, provided that you also convey the machine-readable Corresponding Source under the terms of this License, in one of these ways:

- a.** Convey the object code in, or embodied in, a physical product (including a physical distribution medium), accompanied by the Corresponding Source fixed on a durable physical medium customarily used for software interchange.
- b.** Convey the object code in, or embodied in, a physical product (including a physical distribution medium), accompanied by a written offer, valid for at least three years and valid for as long as you offer spare parts or customer support for that product model, to give anyone who possesses the object code either (1) a copy of the Corresponding Source for all the software in the product that is covered by this License, on a durable

physical medium customarily used for software interchange, for a price no more than your reasonable cost of physically performing this conveying of source, or (2) access to copy the Corresponding Source from a network server at no charge.

- c. Convey individual copies of the object code with a copy of the written offer to provide the Corresponding Source. This alternative is allowed only occasionally and noncommercially, and only if you received the object code with such an offer, in accord with subsection 6b.
- d. Convey the object code by offering access from a designated place (gratis or for a charge), and offer equivalent access to the Corresponding Source in the same way through the same place at no further charge. You need not require recipients to copy the Corresponding Source along with the object code. If the place to copy the object code is a network server, the Corresponding Source may be on a different server (operated by you or a third party) that supports equivalent copying facilities, provided you maintain clear directions next to the object code saying where to find the Corresponding Source. Regardless of what server hosts the Corresponding Source, you remain obligated to ensure that it is available for as long as needed to satisfy these requirements.
- e. Convey the object code using peer-to-peer transmission, provided you inform other peers where the object code and Corresponding Source of the work are being offered to the general public at no charge under subsection 6d.

A separable portion of the object code, whose source code is excluded from the Corresponding Source as a System Library, need not be included in conveying the object code work.

A “User Product” is either (1) a “consumer product”, which means any tangible personal property which is normally used for personal, family, or household purposes, or (2) anything designed or sold for incorporation into a dwelling. In determining whether a product is a consumer product, doubtful cases shall be resolved in favor of coverage. For a particular product received by a particular user, “normally used” refers to a typical or common use of that class of product, regardless of the status of the particular user or of the way in which the particular user actually uses, or expects or is expected to use, the product. A product is a consumer product regardless of whether the product has substantial commercial, industrial or non-consumer uses, unless such uses represent the only significant mode of use of the product.

“Installation Information” for a User Product means any methods, procedures, authorization keys, or other information required to install and execute modified versions of a covered work in that User Product from a modified version of its Corresponding Source. The information must suffice to ensure that the continued functioning of the modified object code is in no case prevented or interfered with solely because modification has been made.

If you convey an object code work under this section in, or with, or specifically for use in, a User Product, and the conveying occurs as part of a transaction in which the right of possession and use of the User Product is transferred to the recipient in perpetuity or for a fixed term (regardless of how the transaction is characterized),

the Corresponding Source conveyed under this section must be accompanied by the Installation Information. But this requirement does not apply if neither you nor any third party retains the ability to install modified object code on the User Product (for example, the work has been installed in ROM).

The requirement to provide Installation Information does not include a requirement to continue to provide support service, warranty, or updates for a work that has been modified or installed by the recipient, or for the User Product in which it has been modified or installed. Access to a network may be denied when the modification itself materially and adversely affects the operation of the network or violates the rules and protocols for communication across the network.

Corresponding Source conveyed, and Installation Information provided, in accord with this section must be in a format that is publicly documented (and with an implementation available to the public in source code form), and must require no special password or key for unpacking, reading or copying.

## 7. ADDITIONAL TERMS.

“Additional permissions” are terms that supplement the terms of this License by making exceptions from one or more of its conditions. Additional permissions that are applicable to the entire Program shall be treated as though they were included in this License, to the extent that they are valid under applicable law. If additional permissions apply only to part of the Program, that part may be used separately under those permissions, but the entire Program remains governed by this License without regard to the additional permissions.

When you convey a copy of a covered work, you may at your option remove any additional permissions from that copy, or from any part of it. (Additional permissions may be written to require their own removal in certain cases when you modify the work.) You may place additional permissions on material, added by you to a covered work, for which you have or can give appropriate copyright permission.

Notwithstanding any other provision of this License, for material you add to a covered work, you may (if authorized by the copyright holders of that material) supplement the terms of this License with terms:

- a.** Disclaiming warranty or limiting liability differently from the terms of sections 15 and 16 of this License; or
- b.** Requiring preservation of specified reasonable legal notices or author attributions in that material or in the Appropriate Legal Notices displayed by works containing it; or
- c.** Prohibiting misrepresentation of the origin of that material, or requiring that modified versions of such material be marked in reasonable ways as different from the original version; or
- d.** Limiting the use for publicity purposes of names of licensors or authors of the material; or

- e. Declining to grant rights under trademark law for use of some trade names, trademarks, or service marks; or
- f. Requiring indemnification of licensors and authors of that material by anyone who conveys the material (or modified versions of it) with contractual assumptions of liability to the recipient, for any liability that these contractual assumptions directly impose on those licensors and authors.

All other non-permissive additional terms are considered “further restrictions” within the meaning of section 10. If the Program as you received it, or any part of it, contains a notice stating that it is governed by this License along with a term that is a further restriction, you may remove that term. If a license document contains a further restriction but permits relicensing or conveying under this License, you may add to a covered work material governed by the terms of that license document, provided that the further restriction does not survive such relicensing or conveying.

If you add terms to a covered work in accord with this section, you must place, in the relevant source files, a statement of the additional terms that apply to those files, or a notice indicating where to find the applicable terms.

Additional terms, permissive or non-permissive, may be stated in the form of a separately written license, or stated as exceptions; the above requirements apply either way.

## 8. TERMINATION.

You may not propagate or modify a covered work except as expressly provided under this License. Any attempt otherwise to propagate or modify it is void, and will automatically terminate your rights under this License (including any patent licenses granted under the third paragraph of section 11).

However, if you cease all violation of this License, then your license from a particular copyright holder is reinstated (a) provisionally, unless and until the copyright holder explicitly and finally terminates your license, and (b) permanently, if the copyright holder fails to notify you of the violation by some reasonable means prior to 60 days after the cessation.

Moreover, your license from a particular copyright holder is reinstated permanently if the copyright holder notifies you of the violation by some reasonable means, this is the first time you have received notice of violation of this License (for any work) from that copyright holder, and you cure the violation prior to 30 days after your receipt of the notice.

Termination of your rights under this section does not terminate the licenses of parties who have received copies or rights from you under this License. If your rights have been terminated and not permanently reinstated, you do not qualify to receive new licenses for the same material under section 10.

## 9. ACCEPTANCE NOT REQUIRED FOR HAVING COPIES.

You are not required to accept this License in order to receive or run a copy of the Program. Ancillary propagation of a covered work occurring solely as a consequence of using peer-to-peer transmission to receive a copy likewise does not require acceptance. However, nothing other than this License grants you permission to propagate or modify any covered work. These actions infringe copyright if you do not accept this License. Therefore, by modifying or propagating a covered work, you indicate your acceptance of this License to do so.

## 10. AUTOMATIC LICENSING OF DOWNSTREAM RECIPIENTS.

Each time you convey a covered work, the recipient automatically receives a license from the original licensors, to run, modify and propagate that work, subject to this License. You are not responsible for enforcing compliance by third parties with this License.

An “entity transaction” is a transaction transferring control of an organization, or substantially all assets of one, or subdividing an organization, or merging organizations. If propagation of a covered work results from an entity transaction, each party to that transaction who receives a copy of the work also receives whatever licenses to the work the party's predecessor in interest had or could give under the previous paragraph, plus a right to possession of the Corresponding Source of the work from the predecessor in interest, if the predecessor has it or can get it with reasonable efforts.

You may not impose any further restrictions on the exercise of the rights granted or affirmed under this License. For example, you may not impose a license fee, royalty, or other charge for exercise of rights granted under this License, and you may not initiate litigation (including a cross-claim or counterclaim in a lawsuit) alleging that any patent claim is infringed by making, using, selling, offering for sale, or importing the Program or any portion of it.

## II. PATENTS.

A “contributor” is a copyright holder who authorizes use under this License of the Program or a work on which the Program is based. The work thus licensed is called the contributor's “contributor version”.

A contributor's “essential patent claims” are all patent claims owned or controlled by the contributor, whether already acquired or hereafter acquired, that would be infringed by some manner, permitted by this License, of making, using, or selling its contributor version, but do not include claims that would be infringed only as a consequence of further modification of the contributor version. For purposes of this definition, “control” includes the right to grant patent sublicenses in a manner consistent with the requirements of this License.

Each contributor grants you a non-exclusive, worldwide, royalty-free patent license under the contributor's essential patent claims, to make, use, sell, offer for sale, import and otherwise run, modify and propagate the contents of its contributor version.

In the following three paragraphs, a “patent license” is any express agreement or commitment, however denominated, not to enforce a patent (such as an express permission to practice a patent or covenant not to sue for patent infringement). To “grant” such a patent license to a party means to make such an agreement or commitment not to enforce a patent against the party.

If you convey a covered work, knowingly relying on a patent license, and the Corresponding Source of the work is not available for anyone to copy, free of charge and under the terms of this License, through a publicly available network server or other readily accessible means, then you must either (1) cause the Corresponding Source to be so available, or (2) arrange to deprive yourself of the benefit of the patent license for this particular work, or (3) arrange, in a manner consistent with the requirements of this License, to extend the patent license to downstream recipients. “Knowingly relying” means you have actual knowledge that, but for the patent license, your conveying the covered work in a country, or your recipient’s use of the covered work in a country, would infringe one or more identifiable patents in that country that you have reason to believe are valid.

If, pursuant to or in connection with a single transaction or arrangement, you convey, or propagate by procuring conveyance of, a covered work, and grant a patent license to some of the parties receiving the covered work authorizing them to use, propagate, modify or convey a specific copy of the covered work, then the patent license you grant is automatically extended to all recipients of the covered work and works based on it.

A patent license is “discriminatory” if it does not include within the scope of its coverage, prohibits the exercise of, or is conditioned on the non-exercise of one or more of the rights that are specifically granted under this License. You may not convey a covered work if you are a party to an arrangement with a third party that is in the business of distributing software, under which you make payment to the third party based on the extent of your activity of conveying the work, and under which the third party grants, to any of the parties who would receive the covered work from you, a discriminatory patent license (a) in connection with copies of the covered work conveyed by you (or copies made from those copies), or (b) primarily for and in connection with specific products or compilations that contain the covered work, unless you entered into that arrangement, or that patent license was granted, prior to 28 March 2007.

Nothing in this License shall be construed as excluding or limiting any implied license or other defenses to infringement that may otherwise be available to you under applicable patent law.

## 12. NO SURRENDER OF OTHERS’ FREEDOM.

If conditions are imposed on you (whether by court order, agreement or otherwise) that contradict the conditions of this License, they do not excuse you from the conditions of this License. If you cannot convey a covered work so as to satisfy simultaneously your obligations under this License and any other pertinent obligations, then as a consequence you may not convey it at all. For example, if you agree to terms that obligate you to collect a royalty for further conveying from those to whom you convey the Program, the only way you could satisfy both those terms and this License would be to refrain entirely from conveying the Program.

### 13. USE WITH THE GNU AFFERO GENERAL PUBLIC LICENSE.

Notwithstanding any other provision of this License, you have permission to link or combine any covered work with a work licensed under version 3 of the GNU Affero General Public License into a single combined work, and to convey the resulting work. The terms of this License will continue to apply to the part which is the covered work, but the special requirements of the GNU Affero General Public License, section 13, concerning interaction through a network will apply to the combination as such.

### 14. REVISED VERSIONS OF THIS LICENSE.

The Free Software Foundation may publish revised and/or new versions of the GNU General Public License from time to time. Such new versions will be similar in spirit to the present version, but may differ in detail to address new problems or concerns.

Each version is given a distinguishing version number. If the Program specifies that a certain numbered version of the GNU General Public License “or any later version” applies to it, you have the option of following the terms and conditions either of that numbered version or of any later version published by the Free Software Foundation. If the Program does not specify a version number of the GNU General Public License, you may choose any version ever published by the Free Software Foundation.

If the Program specifies that a proxy can decide which future versions of the GNU General Public License can be used, that proxy's public statement of acceptance of a version permanently authorizes you to choose that version for the Program.

Later license versions may give you additional or different permissions. However, no additional obligations are imposed on any author or copyright holder as a result of your choosing to follow a later version.

### 15. DISCLAIMER OF WARRANTY.

THERE IS NO WARRANTY FOR THE PROGRAM, TO THE EXTENT PERMITTED BY APPLICABLE LAW. EXCEPT WHEN OTHERWISE STATED IN WRITING THE COPYRIGHT HOLDERS AND/OR OTHER PARTIES PROVIDE THE PROGRAM “AS IS” WITHOUT WARRANTY OF ANY KIND, EITHER EXPRESSED OR IMPLIED, INCLUDING, BUT NOT LIMITED TO, THE IMPLIED WARRANTIES OF MERCHANTABILITY AND FITNESS FOR A PARTICULAR PURPOSE. THE ENTIRE RISK AS TO THE QUALITY AND PERFORMANCE OF THE PROGRAM IS WITH YOU. SHOULD THE PROGRAM PROVE DEFECTIVE, YOU ASSUME THE COST OF ALL NECESSARY SERVICING, REPAIR OR CORRECTION.

## 16. LIMITATION OF LIABILITY.

IN NO EVENT UNLESS REQUIRED BY APPLICABLE LAW OR AGREED TO IN WRITING WILL ANY COPYRIGHT HOLDER, OR ANY OTHER PARTY WHO MODIFIES AND/OR CONVEYS THE PROGRAM AS PERMITTED ABOVE, BE LIABLE TO YOU FOR DAMAGES, INCLUDING ANY GENERAL, SPECIAL, INCIDENTAL OR CONSEQUENTIAL DAMAGES ARISING OUT OF THE USE OR INABILITY TO USE THE PROGRAM (INCLUDING BUT NOT LIMITED TO LOSS OF DATA OR DATA BEING RENDERED INACCURATE OR LOSSES SUSTAINED BY YOU OR THIRD PARTIES OR A FAILURE OF THE PROGRAM TO OPERATE WITH ANY OTHER PROGRAMS), EVEN IF SUCH HOLDER OR OTHER PARTY HAS BEEN ADVISED OF THE POSSIBILITY OF SUCH DAMAGES.

## 17. INTERPRETATION OF SECTIONS 15 AND 16.

If the disclaimer of warranty and limitation of liability provided above cannot be given local legal effect according to their terms, reviewing courts shall apply local law that most closely approximates an absolute waiver of all civil liability in connection with the Program, unless a warranty or assumption of liability accompanies a copy of the Program in return for a fee.

## END OF TERMS AND CONDITIONS

## HOW TO APPLY THESE TERMS TO YOUR NEW PROGRAMS

If you develop a new program, and you want it to be of the greatest possible use to the public, the best way to achieve this is to make it free software which everyone can redistribute and change under these terms.


To do so, attach the following notices to the program. It is safest to attach them to the start of each source file to most effectively state the exclusion of warranty; and each file should have at least the “copyright” line and a pointer to where the full notice is found.

```
one line to give the program's name and a brief idea of what it does.
Copyright (C) year name of author
```

```
This program is free software: you can redistribute it and/or modify
it under the terms of the GNU General Public License as published by
the Free Software Foundation, either version 3 of the License, or
(at your option) any later version.
```



```
This program is distributed in the hope that it will be useful,  
but WITHOUT ANY WARRANTY; without even the implied warranty of  
MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE. See the  
GNU General Public License for more details.
```


You should have received a copy of the GNU General Public License  
along with this program. If not, see [HTTPS://WWW.GNU.ORG/LICENSES/](https://www.gnu.org/licenses/) .


Also add information on how to contact you by electronic and paper mail.

If the program does terminal interaction, make it output a short notice like this when it starts in an interactive mode:

```
program Copyright (C) year name of author  
This program comes with ABSOLUTELY NO WARRANTY; for details type 'show w'.  
This is free software, and you are welcome to redistribute it  
under certain conditions; type 'show c' for details.
```


The hypothetical commands ‘show w’ and ‘show c’ should show the appropriate parts of the General Public License. Of course, your program’s commands might be different; for a GUI interface, you would use an “about box”.

You should also get your employer (if you work as a programmer) or school, if any, to sign a “copyright disclaimer” for the program, if necessary. For more information on this, and how to apply and follow the GNU GPL, see [HTTPS://WWW.GNU.ORG/LICENSES/](https://www.gnu.org/licenses/) .


The GNU General Public License does not permit incorporating your program into proprietary programs. If your program is a subroutine library, you may consider it more useful to permit linking proprietary applications with the library. If this is what you want to do, use the GNU Lesser General Public License instead of this License. But first, please read [HTTPS://WWW.GNU.ORG/LICENSES/WHY-NOT-LGPL.HTML](https://www.gnu.org/licenses/why-not-lgpl.html) .


## COLOPHON

**About the author.** Filippo Rusconi is a senior researcher at the French national research council (*Centre national de la Recherche scientifique*, CNRS). Filippo has a background in biochemistry and organic chemistry and was trained during his Ph.D. as a bioanalytical chemist. He has extensive knowledge of analytical techniques involved in the study of biopolymers.



Filippo Rusconi is the author of a handbook about mass spectrometry for biochemists (French). The book was published by the French sci/tech publisher **LAVOISIER** ([HTTPS://WWW.LAVOISIER.FR](https://www.lavoisier.fr)) .




**Colophon.** The look of this book (PDF file) is the result of me having read many books from the O'Reilly publisher.

The frog on the book title page is a frog from Papua. This frog is able to hover when performing downwards leaps. This picture is courtesy [HTTP://WWW.PAPUAWEB.ORG](http://www.papuaweb.org) .

The typesetting of the book has been done on a Debian GNU/Linux computer using only Free Software. Use of the DocBook Authoring and Publishing Suite (**DAPS** ([HTTPS://GITHUB.COM/OPENSUSE/DAPS](https://github.com/opensuse/daps)) ) from SUSE was key in the process.

The layout adopted for this book is an adaptation of the SUSE stylesheets. I would like to thank Frank Sundermeyer <fsundermeyer@opensuse.org> and Stefan Knorr <sknorr@suse.de> for being helpful with all my questions.

The main font used was **EBGARAMOND** ([HTTPS://GITHUB.COM/GEORGD/EB-GARAMOND](https://github.com/georgd/EB-GARAMOND))  and the symbol/mathematical font was from the **STIX PROJECT** ([HTTPS://WWW.STIXFONTS.ORG/](https://www.stixfonts.org/))  (font: STIX2Math).

The screen shots were taken with Spectacle, the screen capture program shipped along with my **KDE** ([HTTPS://WWW.KDE.ORG/](https://www.kde.org/))  desktop environment and resampled using The GNU image manipulation program **THE GIMP** ([HTTPS://WWW.GIMP.ORG/](https://www.gimp.org/)) . Illustrations were done in **INKSCAPE** ([HTTPS://INKSCAPE.ORG/](https://inkscape.org/)) , a vectorial drawing software.