**Intro Guide to Analyzing Metabolomic Data using mzMine 2.53 and GNPS**

**MichiganState 2020 iGEM | Bioinformatics**

**Last Updated: 10/24/2020**

**Preface**

This short document is meant to share our feature-based metabolomic networking workflow/protocol (from ingesting raw .mzML mass spec files to exporting data for analysis in GNPS) and does not include any statistical processes to validate our results. This guide is primarily targeted towards new users of mzMine/GNPS; and as such, we've included our raw data on our wiki so that you can follow this protocol. Furthermore, we assume readers understand the basic concepts of molecular networking prior to opening this doc; but just in case, we've linked more resources below that offer more comprehensive explanations of FBMN methods. Additionally, the parameters used for this analysis may not be appropriate for all situations, and we cannot guarantee when the information contained in this document will be deprecated. In anycase, we hope the information enclosed in this document will be informative to future iGEM teams looking to incorporate metabolomics analysis to their projects.

**MZmine/GNPS FBMN Outline/Table of Contents/Summary w/specific parameters used**
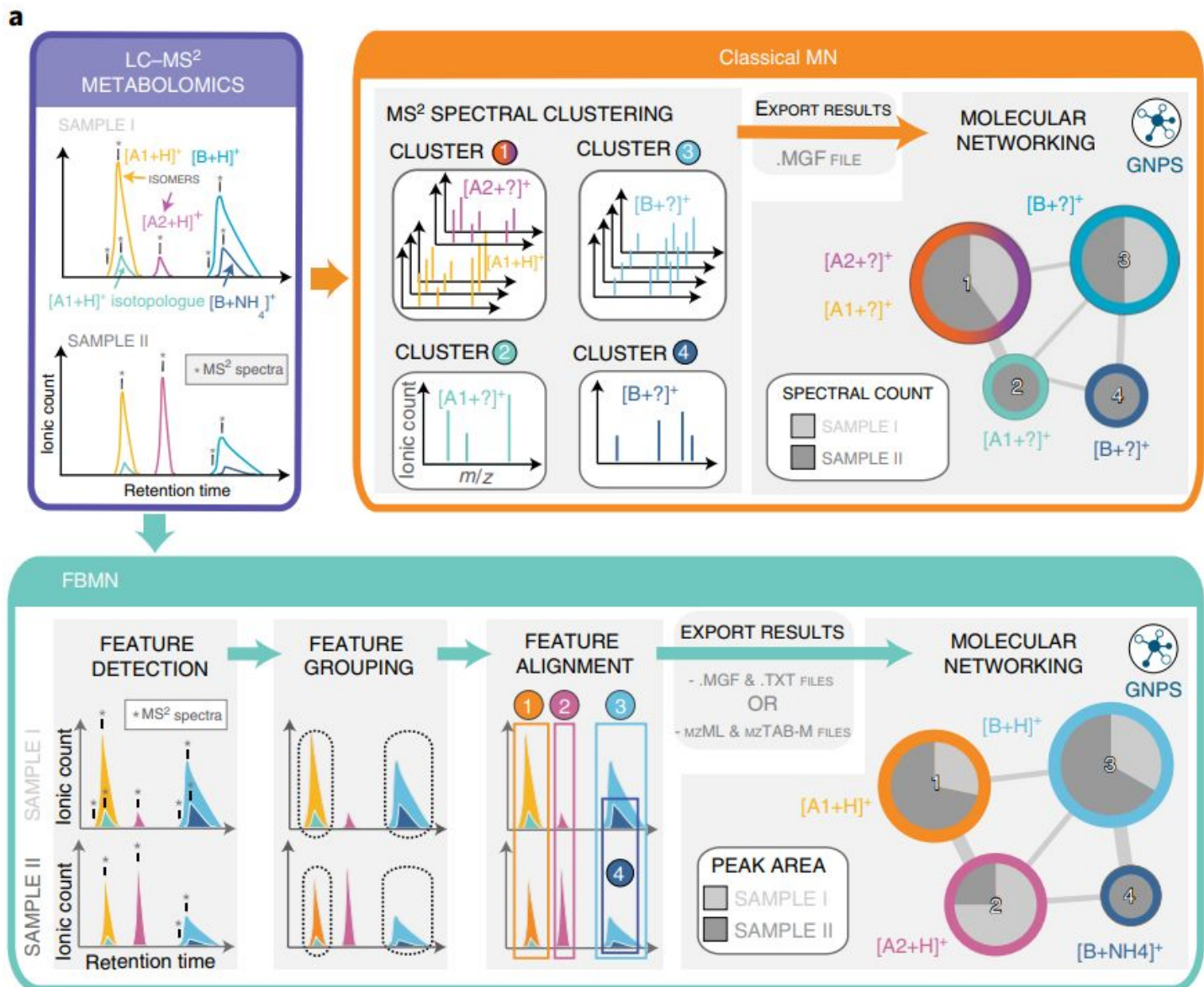
## Graphical Overview



Figure 1A. Graphical overview of FBMN compared to Classical MN. Adapted [reprinted] from "Feature-based molecular networking in the GNPS analysis environment," by Nothias et. al., 2020, *Nature Methods,17,* pp. 905-908.

**Prerequisites**

**Learn about Feature-Based Molecular Networking**

- This paper by Nothias et al. (2020) offers a great comparison between FBMN and classical networking and offers insight into how you can use FBMN within the GNPS Environment.
  - Nothias, L., Petras, D., Schmid, R. et al. Feature-based molecular networking in the GNPS analysis environment. Nat Methods 17, 905–908 (2020). https://doi.org/10.1038/s41592-020-0933-6
- This more extensive tutorial is a great follow-up example that goes more in-depth in more complex topics like manual validation of feature lists generated in the MZmine workflow.
  - **https://ccms-ucsd.github.io/GNPSDocumentation/featurebasedmolecularnetworking-with-mzmine2/**

**Get mzMine**

- Follow instruction here to download the latest version of MZmine: https://github.com/mzmine/mzmine2/releases
  - Installation isn't necessary as you initialize the application from start-up command/batch files for your respective operating system.
    - [Please note the images on this doc are generated on Windows, but the GUI does not change between operating systems]
  - More in depth documentation on starting and operating MZmine can be found within this manual:
    - https://docs.google.com/document/d/1JtdNwp1y6wz_Boz8qVGESH-kSrMhcGaJWR18iGmhstU/edit?usp=sharing
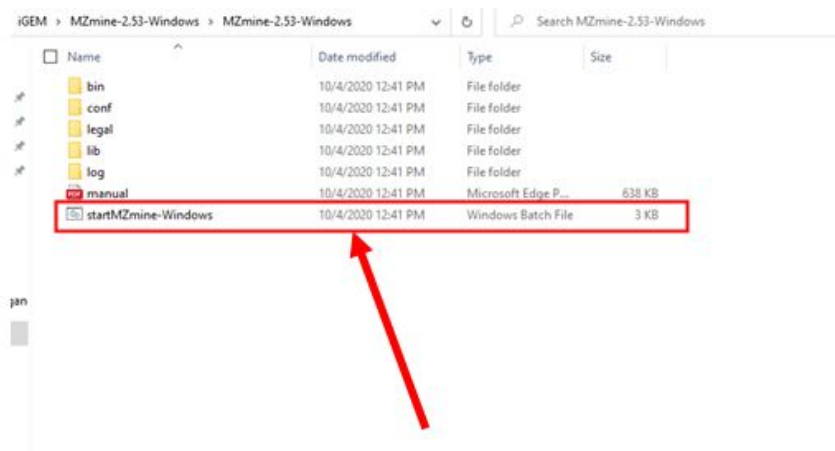
**Set-up Account on GNPS**

- https://gnps.ucsd.edu/ProteoSAFe/static/gnps-splash.jsp?redirect=auth
- More details about GNPS can be found in its Documentation page:
  - https://ccms-ucsd.github.io/GNPSDocumentation/

**Setup file transfer client to upload data to GNPS**

- Use these instructions to easily upload files to your GNPS account to analyze your samples.
  - https://ccms-ucsd.github.io/GNPSDocumentation/fileupload/
- FileSharing Clients make these file transfers easier:
  - Filezilla: https://filezilla-project.org/
  - WinSCP: https://winscp.net/eng/download.php

## Ingest Data to MZmine

- After unzipping the mzMine program, initialize the program by clicking on the startMZmine script:



- Once initialized, you should see the following window:

**Import Raw Data Files**

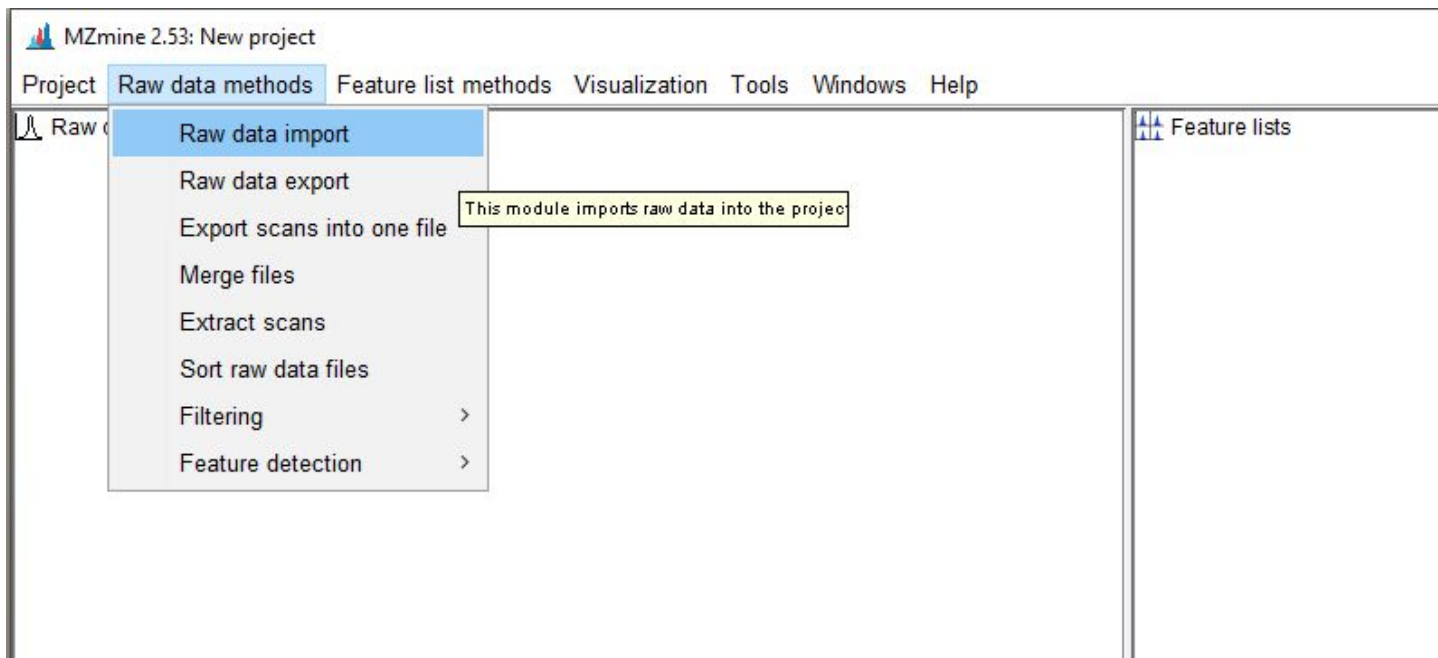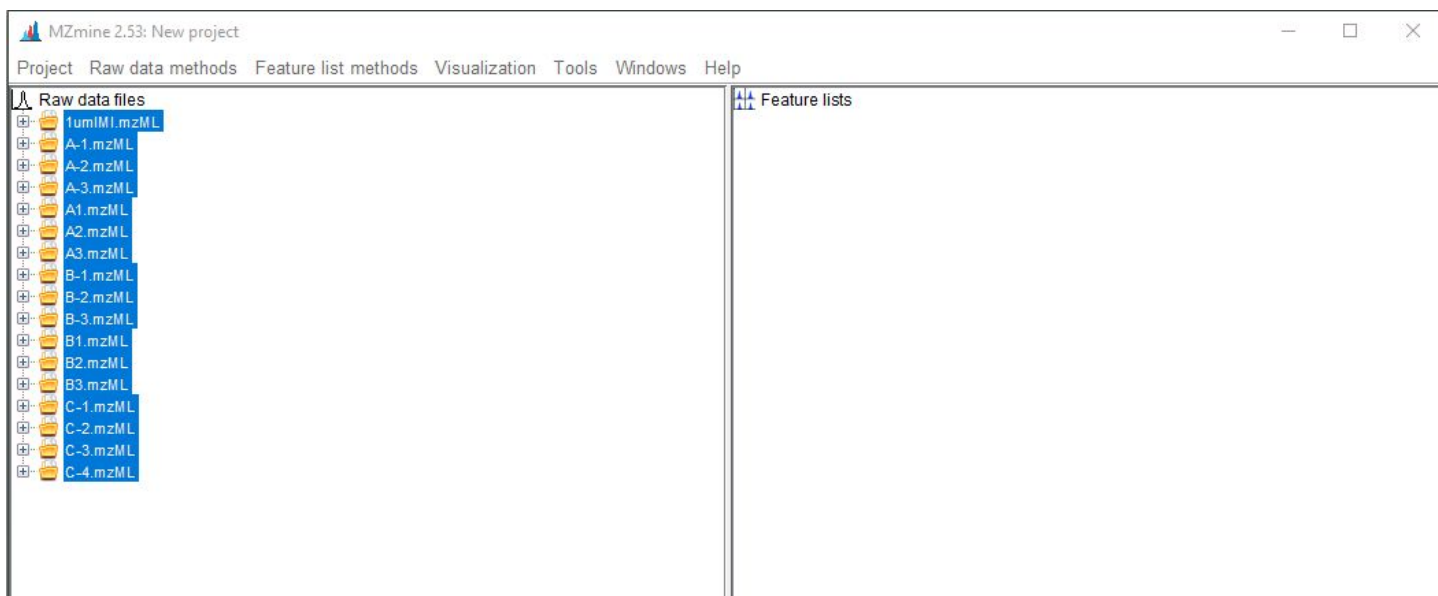- Start importing your .mzXML files by clicking the "Raw data import" tool under "Raw Data methods" tab and selecting your files from your files:
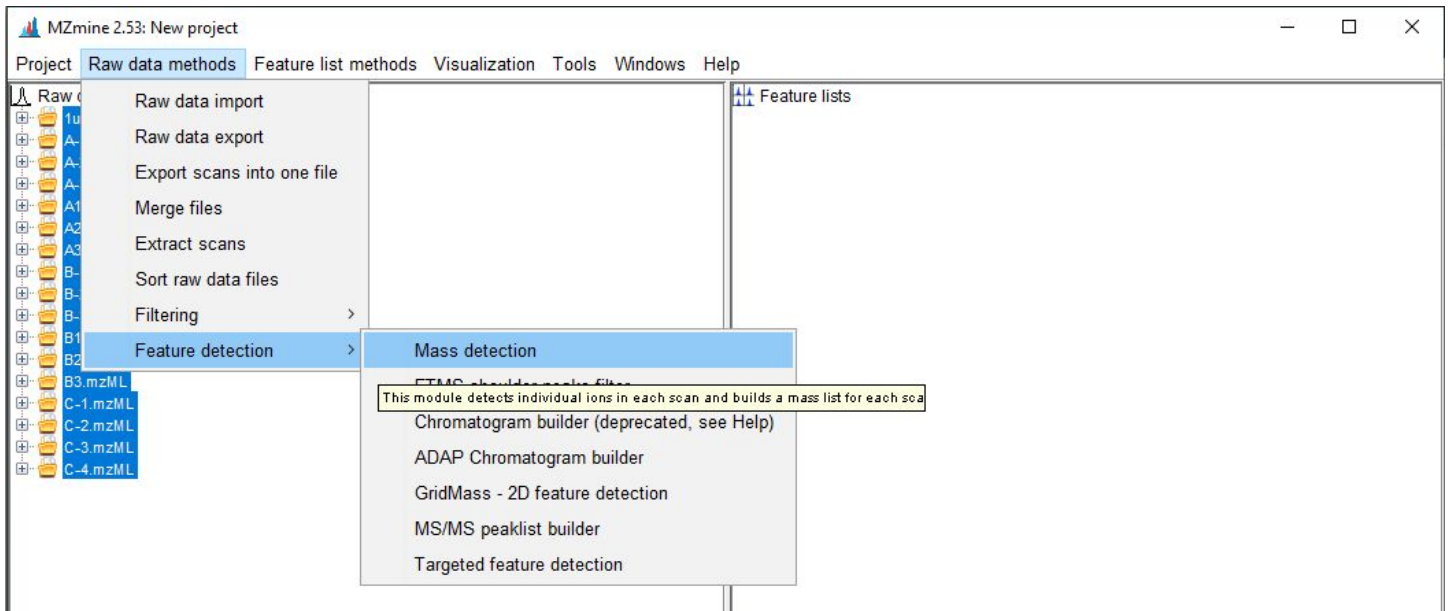


- Once, imported your raw data should show up in the left panel as individual folders that include the individual Mass-Spec files:

## Mass Detection

- Select all of of the raw data files and click on the "Mass Detection" tool using the following path: "Raw data methods" → "Feature Detection" → "Mass Detection":



- The following pop-up box should appear:
    - Under Raw data files, select the files you want to analyze that were originally selected in the main window.
    - Mass Detector should be set to "Centroid"
    - Mass list name should be "masses"
        - This step filters the chromatograms based on a specific noise threshold and creates folders that are named after the original raw data files with "masses" appended to the end of the selected files

- Set the MS1 noise levels by clicking on the three dots to the right of the dropbox menu of the Mass Detector field:



- The following pop-up box should appear (note that the "Show preview" option is clicked):
  - Set the noise level appropriate for the run, in this case we set the intensity threshold to 1E4 (1 x 10$^6$)
  - In the preview box, you can see that the peaks that do not reach this threshold are not highlighted in red.

- Repeat the last two steps, but now filtering the MS 2.
  - Change from setting noise levels from MS 1 to MS 2 by clicking the "Set Filter" box under the "Scans" field:



  - The noise level should be set such that most of the noise is filtered out of your samples, but (in this example) it is not a quantitative step.
    - In general though, the MS2 scan level should be set less than the MS1 level.
    - For our files, we set noise levels to **1e4 and 1e3 for the MS1 and MS2** levels, respectively.

* Note how the file icons change once you apply the filters with a green check mark on each file icon.

**ADAP Chromatogram Builder**

- Next, we will build the chromatograms from our mass spec files:
  - Clicking on the ADAP Chromatogram builder from the following path: "Raw data methods" →
    "Feature detection" → "ADAP Chromatogram Builder":



- For our experiment we set the following parameters to build our chromatograms:
  - Select the mass list by clicking the "Choose…" button to the right of the mass list option and picking
    the "masses" option in the dropdown menu.
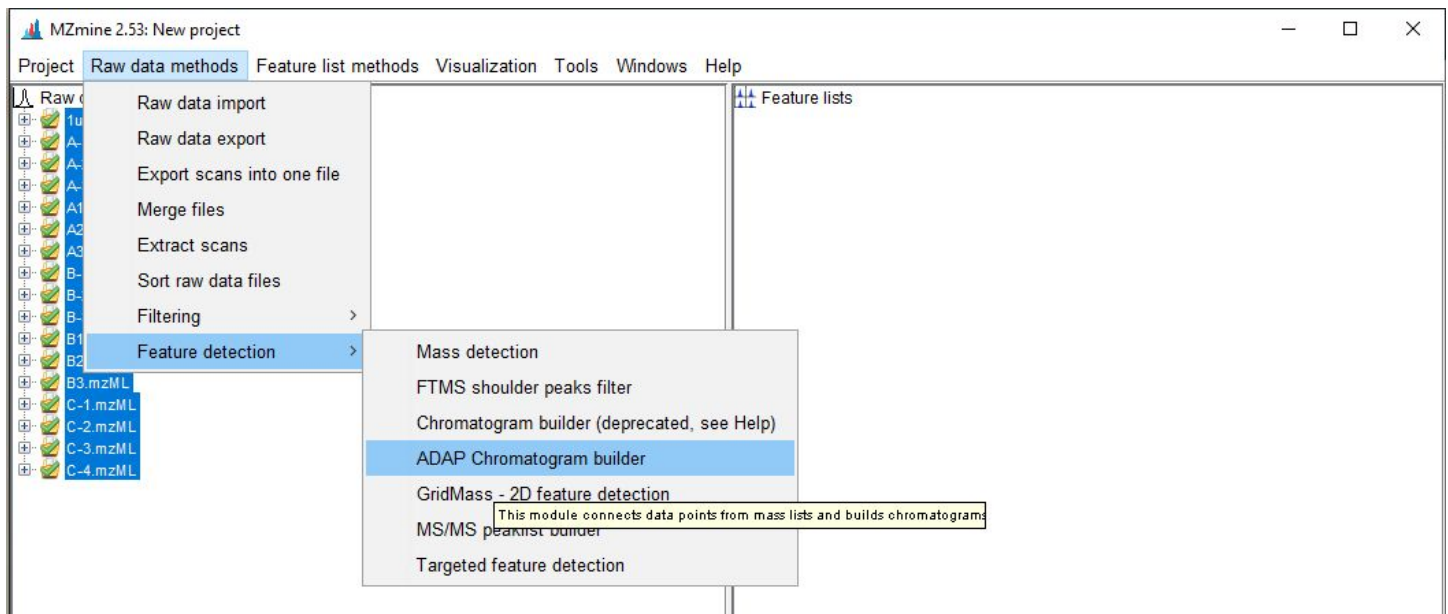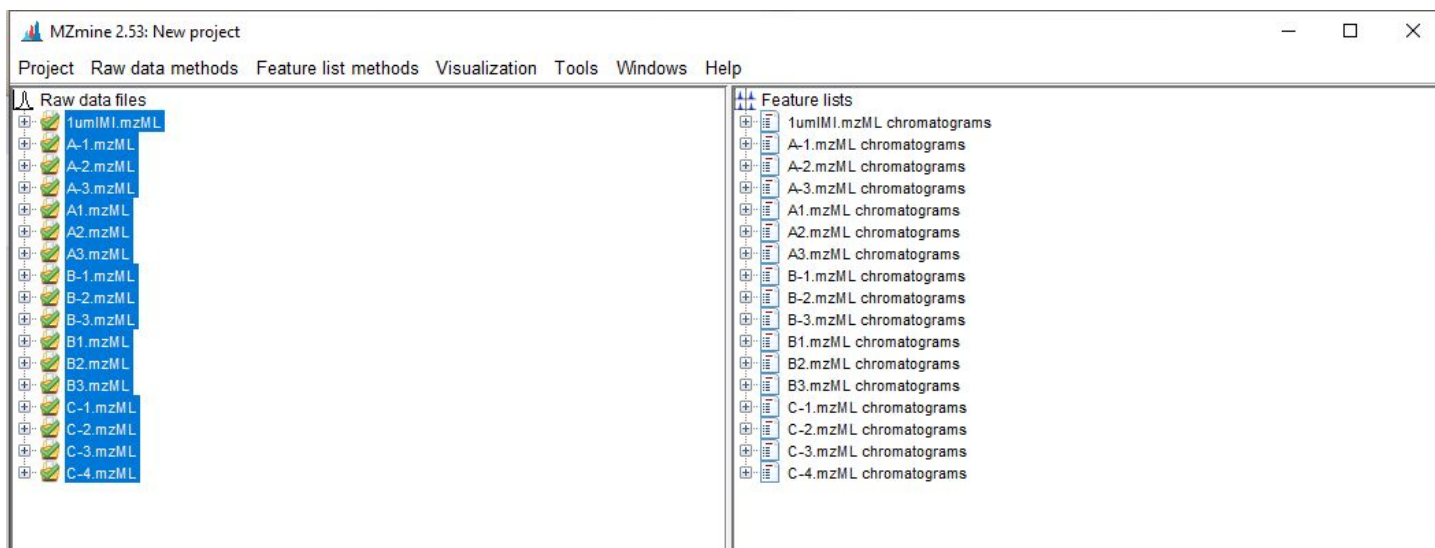- More details about these parameters can be found in this manual
  (http://mzmine.github.io/ADAP_user_manual.pdf), but the following explanations are excerpted briefly:
  - Min group size in # of Scans: 5
    - There must be this number of sequential scans having points above a set intensity threshold.
  - Group Intensity Threshold: 3e6
    - The threshold by which scans are compared against to count spectra as part of a group
  - Min highest intensity: 3e6
    - The group of spectra must have a peak of this intensity to be counted
  - m/z tolerance: 0.02 m/z or 10 ppm
    - Tolerance of the difference between data points in consecutive scans to be counted
  - Suffix: chromatograms
    - Appends the file name with "chromatograms"
- After entering these parameters, press ok to build your chromatograms.

- The chromatograms will appear on the right in the "Feature lists" panel:



**Chromatographic Deconvolution**
- Next, step is Chromatogram deconvolution
  - Chromatogram Deconvolution is the process of separating various mass spec data from the measured chromatographic information. Oftentimes in LC, there are samples that will co-elute in the same time period thereby "convoluting" the resulting mass spec data. Deconvolution reverses this process to separate the various masses that are co-eluted at the same time point.
    - Check out this page for a more in depth explanation of this process:
      https://www.spectralworks.com/products/analyzerpro-xd/learn-about-analyzerpro-xd/chromatographic-deconvolution/
- The Chromatogram Deconvolution Tool can be found under the following path: "Feature list methods" → "Feature Detection" → "Chromatogram deconvolution"
  - Be sure to select the chromatograms before opening this tool

- You should see the following pop-out



- We will use the following parameters:
  - Suffix: deconvoluted
  - Algorithm: Baseline cut-off (*Click the three dots next to the Baseline cut-off algorithm to modify the Algorithm parameters)
    - Min Peak height: 1e5
    - Peak Duration: 0 to 3 minutes
    - Baseline Level: 1e4
    - In the preview window, you can see individual peaks are color coded and you can adjust your parameters to ensure most of the features are identified and deconvoluted.
    - Individual peaks are identified by the m/z (mass to charge ratio) and retention time ( in minutes)

- ○ m/z center calculation: MEDIAN
- ○ All others: leave blank

**Isotope Grouping**
- ● Next, we will identify groups of isotopes with the "Isotopic peak grouper" tool which can be found under the following path: "Feature list methods" → "Isotopes" → "Isotopic peak grouper".

- The following pop-up will appear, and we'll use the following parameters to group isotopes:
    - Name suffix: deisotoped
    - m/z tolerance: 0.02 or 10 ppm
    - RT tolerance: 0.25 (absolute (min))
    - max charge: 2
    - Representative isotope: Most intense

| Please set the parameters | | ✕ |
|---|---|---|
| Feature lists | As selected in main window | ... |
| Name suffix | deisotoped | |
| m/z tolerance | 0.02   m/z  or  10 | ppm |
| Retention time tolerance | 0.25   absolute (min) | |
| Monotonic shape | ☐ | |
| Maximum charge | 2 | |
| Representative isotope | Most intense | |
| Remove original peaklist | ☐ | |

OK    Cancel    Help

**Join Aligner**

- Next, we will align these groups of isotopes and other detected features according to their mass spec peaks and retention time for each sample.
  - Classical molecular networking primarily relies on the alignment of mass spec peaks to develop relationships between metabolites (as seen in the graphical overview).
- The alignment tools can be found with the following path: "feature list methods" → "Alignment"
  - For our samples we will use the simple Join aligner



- The following Pop-up should appear:



- And we'll use the following parameters:
  - Feature list name: Aligned feature list
  - m/z tolerance: 0.02 m/z or 10 ppm
  - Weight for m/z: 75

- - ■ The is the weight of mass spec information on the scoring function fo the join aligner.
  - ○ Retention time tolerance: 0.25 [absolute (min)]
    - ■ This is the maximum allowed difference in LC-peak retention time between matched species.
  - ○ Weight for RT: 25
    - ■ Like above, this refers to the weight of RT on the scoring function.
- The product of this step result in a single list "Aligned feature list" which can be viewed in a separate window:
  - ○ Your various samples will be displayed along the top, and the individual identified peaks will be arranged according to an "ID" number assigned during our Feature list method steps.
  - ○ Some interesting quantities you should watch out for:
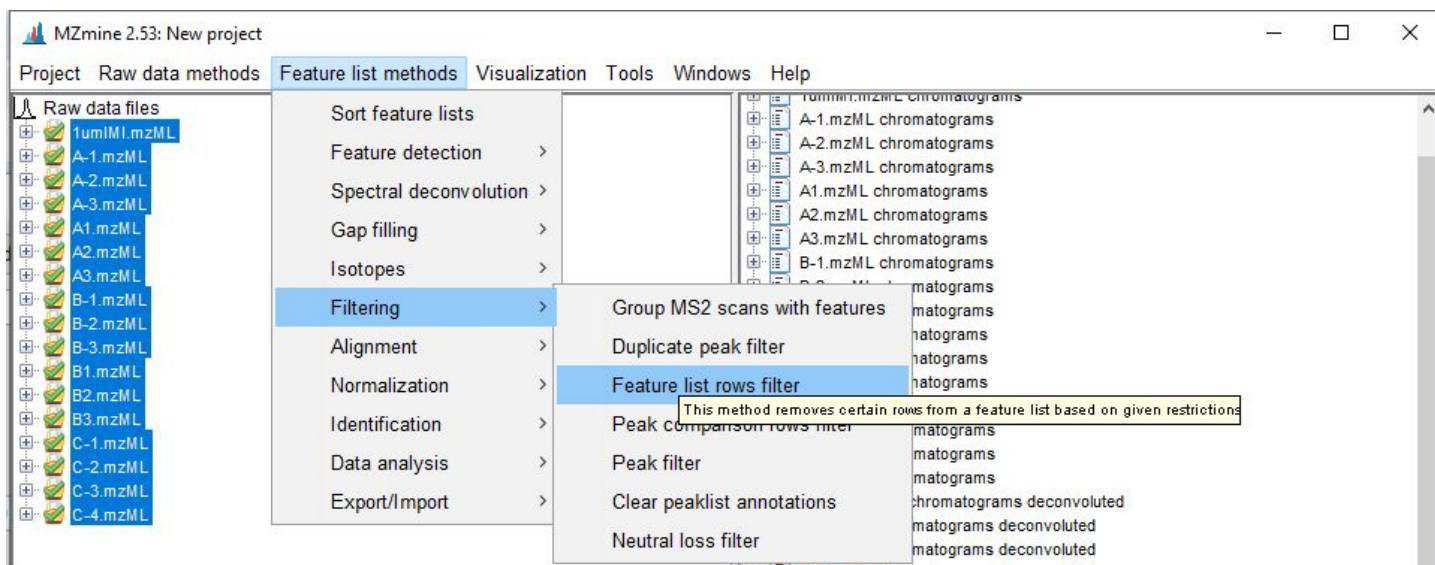    - ■ m/z ratio, retention time (RT), Height, and Area
  - ○ For example in the following peak ID (to save space, only one sample from triplicates of the negative control and treated samples are shown below):
    - ■ The features grouped in this node corresponding to a unique species have an average m/z of ~511and a retention time of ~4 minutes.
    - ■ The height, or intensity of chromatographic peak, is 2.9E7 for our imidacloprid (IMI) standard and 6.0E8 for our *Pseudomonas* sample exposed to imidacloprid.
    - ■ And the area of the peak under the chromatographic peak is similarly larger for our sample exposed to imidacloprid where our concentration of Imidacloprid is higher (1.76 mM vs 1 uM IMI)
    - ■ Using this relationship, we designed future experiments to semi-quantitatively characterize the degradation of Imidacloprid.

| ID | Average | | | Peak shape | 1umIMI.mzML | | | A-1.mzML | | | A1.mzML | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | m/z | RT | | | Status | Height | Area | Status | Height | Area | Status | Height | Area |
| 22 | 511.1128 | 4.37 | | | 🟢 | 2.9E7 | 9.7E7 | 🔴 | | | 🟢 | 6.0E8 | 4.4E9 |

**Peak Filtering**
- The previous aligned feature list step identified over 22,000 nodes/species, of which many are not necessary for downstream analysis. This step allows users to selectively choose which species to analyze in later analysis steps. For our purposes, we will be filtering species based on the number of peaks for each feature, and the m/z of the species. For our project, we have relied on visualizing networks within the GNPS environment, and this step also allows us to introduce several processing steps to make it easier to upload data to GNPS.
- The peak filtering method for this example, "Feature list rows filter", can be found through the following path: "Feature list methods" → "Filtering" → "Feature list rows filter"

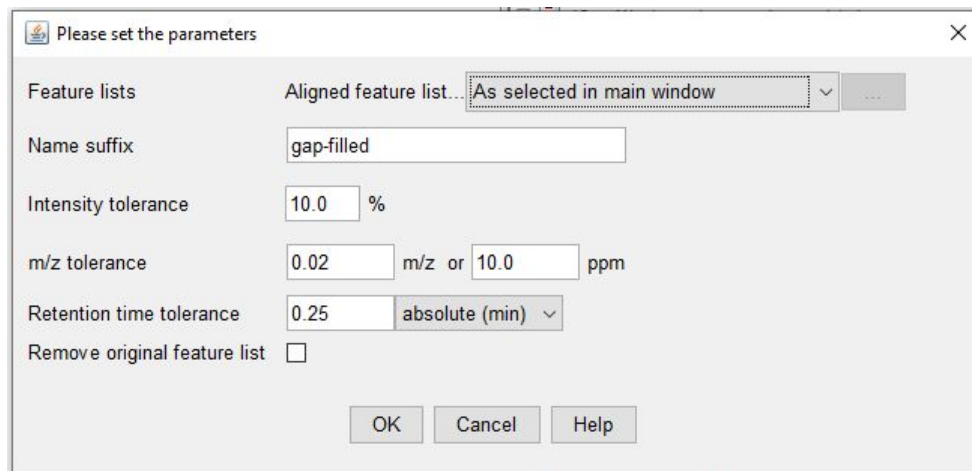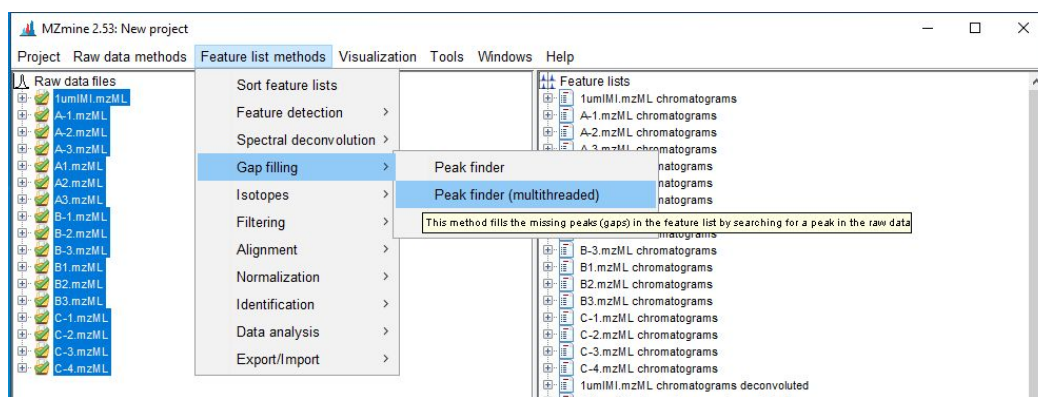- The following pop-up will appear, and we'll use the following parameters (leaving all other default settings unchanged):
  - min peak in row: 5
  - m/z range: 75-1500
  - keep peaks with MS2 for GNPS: TRUE
  - reset peak number ID: TRUE

**Gap-filling:**
- The gap-filling algorithm ("Peak finder") searches through the original raw data to find spectra that may have been mis-aligned during the alignment step or spectra that were filtered out due to low peak intensity (or a previous algorithm may have accidentally removed an important identifying peak).
- For our samples, we'll use the following parameters, and the gap-filling tool can be found by the following path: "Feature list methods" → "Gap filling" → "Peak finder"
  - *The Peak finder tool and the multithreaded peak finder tool use the same parameters and are interchangeable, but the speed of this step can be sped up by allocated tasks to more CPU cores with the multithreaded tool.
  - Intensity tolerance: 10%
  - m/z tolerance: 0.02 or 10 ppm
  - RT tolerance: 0.25



- Visually inspecting the resulting list "Aligned feature list filtered gap-filled" shows which peaks are "estimated" with a yellow status indicator vs the typical green and red status indicator (depending on your gap-filling parameters) as shown in the original aligned feature list.
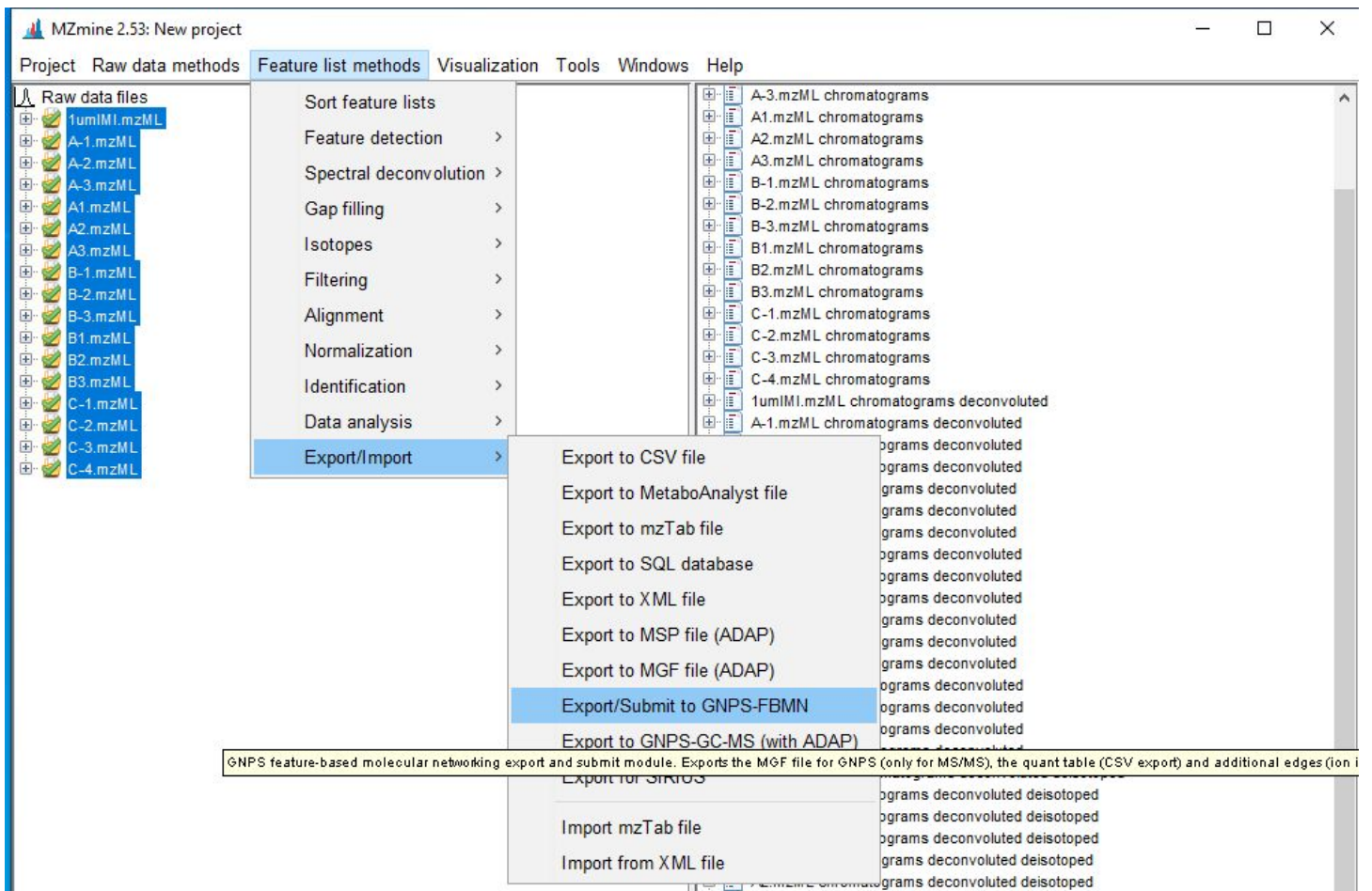
**\*Manual Validation**
- Unfortunately, this topic is beyond the scope of this introductory guide.
- In the simplest terms, you should manually inspect the peaks (as shown under "peak" in feature lists to ensure there are no major abnormalities in the peaks identified (e.g. a broad, undefined peak vs. a distinct peaks

**Export and Uploading Data to GNPS:**

- The team behind the development of GNPS have great documentation on this part, but we briefly outline the steps here for documentation purposes for our team.
    - This video on YouTube walks through the process of uploading data to GNPS after processing data on MZmine: https://youtu.be/vFcGG7T_44E
    - This documentation page walks through another example project and walks through the process of FBMN workflow: https://ccms-ucsd.github.io/GNPSDocumentation/featurebasedmolecularnetworking-with-mzmine2/
        - The process outlined in this page is largely similar, with the exception of the addition of Gap filling and MS row filtering found in our workflow.

**Export:**

- MZmine has several convenient exporting tools that will output data in specific file formats or even packages of data to be analyzed in external tools like GNPS and SIRIUS.
- For our purposes we'll use the GNPS option which can be found through the following path: "Feature list methods" → "Export/Import" → "Export/Submit to GNPS-FBMN"



- The following pop-up should appear; and for our example, we'll use the following parameters:
  - Feature lists: as selected
  - Mass list: masses (*click "Choose…" button like in previous steps and select masses")
  - Filter rows: ALL (you can export a smaller file and just export the MS2 spectra, but including MS1 may improve your odds of features being identified in GNPS)

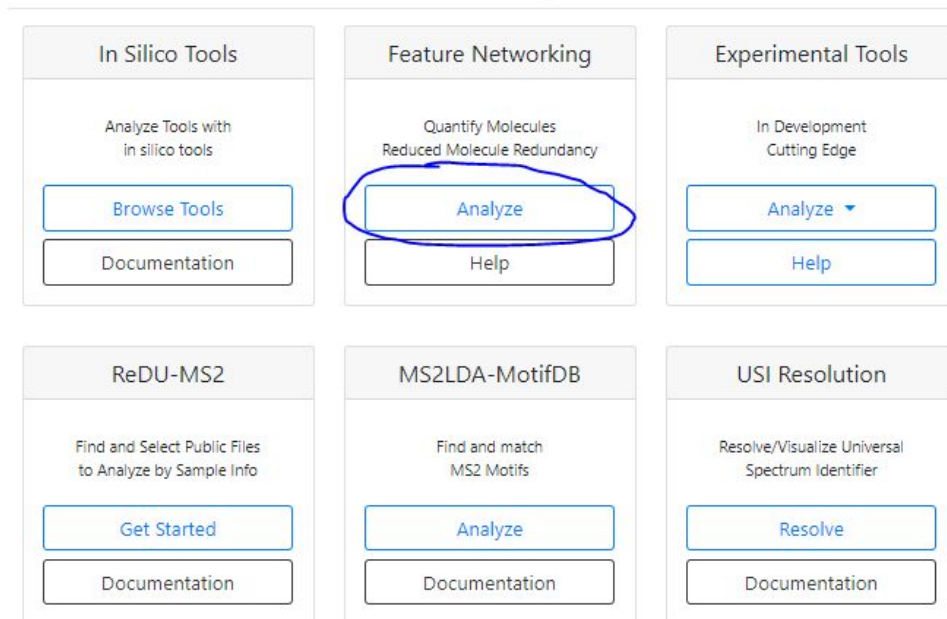- The two output files required for GNPS-FBMN will appear in the folder in which you un-zipped MZmine:



**Uploading files to GNPS:**
- Using the file transfer client of choice, upload the .mgf data and _quant.csv quantification table to your personal storage account
- On GNPS's main home page (after logging into your account in your web-browser), fin "Feature Networking" tool under "Advanced Analysis Tools":

# Advanced Analysis Tools



- Give your job a title (under "Workflow Selection")
- Under "File Selection", when you click "Select Input Files", you can select input files front he files you uploaded in the first step of this section with the following pop-up browser



- After selecting your data to run in the job, you can customize a lot of other parameters, but for this example, we will leave the remainder default values.

- After pressing, submit, you should see the following status page:

| Job Status | |
|---|---|
| **Workflow** | FEATURE-BASED-MOLECULAR-NETWORKING (version release_26) |
| **Status** | RUNNING<br>[Clone] [Clone to Latest Version]                                                    [Delete] |
| **User** | pascualn (pascualn@msu.edu), Michigan State University |
| **Title** | FBMN Example (First Trial Data)wor |
| **Date Created** | 2020-10-23 17:40:37.0 |
| **Execution Time** | 16 seconds |
| **Progress** | |
| **MS2 File MGF/MSP(Progenesis QI)/mzML(MzTab-M)** | pascualn/iGEM/FirstTrial/FBMN_Pascual_1/FBMN_Pascual_firsttry.mgf |
| **Feature Quantification Table** | pascualn/iGEM/FirstTrial/FBMN_Pascual_1/FBMN_Pascual_firsttry_quant.csv |
| **Original mzML Files** | pascualn/iGEM/FirstTrial/AllmzML/1umIMI.mzML<br>pascualn/iGEM/FirstTrial/AllmzML/C-4.mzML<br>pascualn/iGEM/FirstTrial/AllmzML/A-1.mzML<br>pascualn/iGEM/FirstTrial/AllmzML/A-2.mzML<br>pascualn/iGEM/FirstTrial/AllmzML/A-3.mzML<br>pascualn/iGEM/FirstTrial/AllmzML/A1.mzML<br>pascualn/iGEM/FirstTrial/AllmzML/A2.mzML<br>pascualn/iGEM/FirstTrial/AllmzML/A3.mzML<br>pascualn/iGEM/FirstTrial/AllmzML/B-1.mzML<br>pascualn/iGEM/FirstTrial/AllmzML/B-2.mzML<br>pascualn/iGEM/FirstTrial/AllmzML/B-3.mzML<br>pascualn/iGEM/FirstTrial/AllmzML/B1.mzML<br>pascualn/iGEM/FirstTrial/AllmzML/B2.mzML<br>pascualn/iGEM/FirstTrial/AllmzML/B3.mzML<br>pascualn/iGEM/FirstTrial/AllmzML/C-1.mzML<br>pascualn/iGEM/FirstTrial/AllmzML/C-2.mzML<br>pascualn/iGEM/FirstTrial/AllmzML/C-3.mzML |
| **Sample Metadata Table** | pascualn/iGEM/FirstTrial/Metadata/AllMetadata.tsv |
| **Spectral Library** | speclibs/MIADB/MIADB.mgf<br>speclibs/BILELIB19/BILELIB19.mgf<br>speclibs/GNPS-NIH-NATURALPRODUCTSLIBRARY_ROUND2_POSITIVE/GNPS-NIH-NATURALPRODUCTSLIBRARY_ROUND2_POSITIVE.mgf<br>speclibs/GNPS-COLLECTIONS-PESTICIDES-NEGATIVE/GNPS-COLLECTIONS-PESTICIDES-NEGATIVE.mgf<br>speclibs/MASSBANK/MASSBANK.mgf<br>speclibs/UM-NPDC/UM-NPDC.mgf<br>speclibs/GNPS-NIH-CLINICALCOLLECTION1/GNPS-NIH-CLINICALCOLLECTION1.mgf<br>speclibs/HMDB/HMDB.mgf<br>speclibs/DEREPLICATOR_IDENTIFIED_LIBRARY/DEREPLICATOR_IDENTIFIED_LIBRARY.mgf<br>speclibs/GNPS-SELLECKCHEM-FDA-PART1/GNPS-SELLECKCHEM-FDA-PART1.mgf<br>speclibs/RESPECT/RESPECT.mgf<br>speclibs/GNPS-FAULKNERLEGACY/GNPS-FAULKNERLEGACY.mgf<br>speclibs/MMV_NEGATIVE/MMV_NEGATIVE.mgf<br>speclibs/MASSBANKEU/MASSBANKEU.mgf<br>speclibs/GNPS-LIBRARY/GNPS-LIBRARY.mgf<br>speclibs/GNPS-PRESTWICKPHYTOCHEM/GNPS-PRESTWICKPHYTOCHEM.mgf<br>speclibs/GNPS-MSMLS/GNPS-MSMLS.mgf<br>speclibs/GNPS-EMBL-MCF/GNPS-EMBL-MCF.mgf<br>speclibs/PSU-MSMLS/PSU-MSMLS.mgf<br>speclibs/GNPS-IOBA-NHC/GNPS-IOBA-NHC.mgf<br><br>More available |

- And once complete, you can browse through the various GNPS outputs through this link:
    - https://gnps.ucsd.edu/ProteoSAFe/status.jsp?task=85753400a82046728290c4190b850bed

## References:

Katajamaa, M., Miettinen, J. & Oresic, M. MZmine: toolbox for processing and visualization of mass spectrometry based molecular profile data. Bioinformatics 22, 634–636 (2006).

Myers OD, Sumner SJ, Li S, Barnes S, Du X: One Step Forward for Reducing False Positive and False Negative Compound Identifications from Mass Spectrometry Metabolomics Data: New Algorithms for Constructing Extracted Ion Chromatograms and Detecting Chromatographic Peaks. Anal Chem 2017, DOI: 10.1021/acs.analchem.7b00947

Nothias, L., Petras, D., Schmid, R. et al. Feature-based molecular networking in the GNPS analysis environment. Nat Methods 17, 905–908 (2020). https://doi.org/10.1038/s41592-020-0933-6

Pluskal, T., Castillo, S., Villar-Briones, A. & Oresic, M. MZmine 2: modular framework for processing, visualizing, and analyzing mass spectrometry-based molecular profile data. BMC Bioinformatics 11, 395 (2010).

Ty N.F. Roach, Jenna Dilworth, H. Christian Martin, A. Daniel Jones, Robert Quinn, Crawford Drury. bioRxiv 2020.05.10.087072; doi: https://doi.org/10.1101/2020.05.10.087072

Wang, M. et al. Sharing and community curation of mass spectrometry data with Global Natural Products Social Molecular Networking. Nat. Biotechnol. 34, 828–837 (2016).