# Tutorial of validating KGMN unknowns with repository mining

Zhiwei Zhou, Ph.D., 2022.06.08

This tutorial aims to help users to select and validate their interesting unknown peaks from KGMN through repository mining. In the manuscript, we mainly used **MASST** to perform repository mining. The MASST[1] is a tool to query spectrum in context of where it occurs against all GNPS data sets. In this tutorial, we focus on demonstrating how to combine KGMN results and MASST. The detail instructions of MASST can be found in **GNPS document**.

The step-by-step instruction has been provided below.

## 1. Data preparing.

In this workflow, the data files require KGMN (MetDNA2) processed firstly. Here, we utilized NIST human urine data set as example. The data set has been analyzed with KGMN (v1.0.4), and the results can be downloaded **here**.

The folders should look like as below:

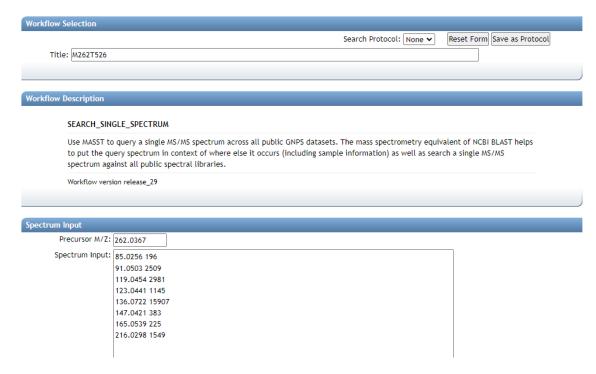| Name | Date modified | Type | Size |
| --- | --- | --- | --- |
| 00_annotation_table | 6/6/2022 2:54 PM | File folder | |
| 02_result_MRN_annotation | 6/6/2022 2:54 PM | File folder | |
| 04_biology_intepretation | 6/4/2022 3:36 PM | File folder | |
| 05_analysis_report | 6/6/2022 2:54 PM | File folder | |
| 06_visualization | 6/6/2022 2:54 PM | File folder | |
| data.csv | 1/17/2022 9:12 AM | Microsoft Excel C... | 2,385 KB |
| NIST_urine01_pos-NIST_urine01.mgf | 1/17/2022 9:10 AM | MGF File | 9,877 KB |
| NIST_urine02_pos-NIST_urine02.mgf | 1/17/2022 9:12 AM | MGF File | 9,895 KB |
| NIST_urine03_pos-NIST_urine03.mgf | 1/17/2022 9:12 AM | MGF File | 9,921 KB |
| NIST_urine04_pos-NIST_urine04.mgf | 1/17/2022 9:10 AM | MGF File | 9,936 KB |
| para_list.txt | 6/4/2022 3:33 PM | Text Document | 2 KB |
| QC_pos-QC.mgf | 1/17/2022 9:12 AM | MGF File | 9,687 KB |
| RT_recalibration_table.csv | 1/17/2022 9:12 AM | Microsoft Excel C... | 1 KB |
| sample.info.csv | 1/17/2022 9:12 AM | Microsoft Excel C... | 1 KB |

The users can browser and select interesting known/unknown peaks in the **annotation table "table1_identification.csv"** in the "00_annotation_table" folder. It should be note that the selection of targeted peak is customized.

For demonstration, we utilized the unknown peak M262T526 as an example (Figure 5d in manuscript). The MS/MS spectrum of this peak can be found in the **"ms2_data.msp"** in "06_visualization" folder. You can open it with text tool (e.g. Notepad++).
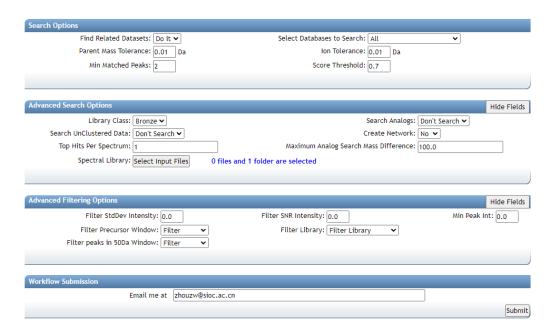
```
7925    NAME: M262T526
7926    PRECURSORMZ: 262.0367
7927    IONMODE: positive
7928    RETENTIONTIME: 526.026
7929    Links:
7930    Comment:
7931    Num Peaks: 8
7932    85.0256 196
7933    91.0503 2509
7934    119.0454 2981
7935    123.0441 1145
7936    136.0722 15907
7937    147.0421 383
7938    165.0539 225
7939    216.0298 1549
```

## 2. Upload and analysis in MASST.

Users can upload this file to MASST (https://gnps.ucsd.edu/ProteoSAFe/static/gnps-splash.jsp?redirect=auth) to perform repository mining. The users need to login first. Then, click the "**query spectrum**" button in MASST panel to start the analysis. Copy **related texts from MSP** file to "title", "precursor m/z", "spectrum input" panel in the web server, respectively.

**Workflow Selection**

Search Protocol: None ▾  Reset Form  Save as Protocol

Title: M262T526

**Workflow Description**

SEARCH_SINGLE_SPECTRUM

Use MASST to query a single MS/MS spectrum across all public GNPS datasets. The mass spectrometry equivalent of NCBI BLAST helps to put the query spectrum in context of where else it occurs (including sample information) as well as search a single MS/MS spectrum against all public spectral libraries.

Workflow version release_29

**Spectrum Input**

Precursor M/Z: 262.0367

Spectrum Input:
```
85.0256 196
91.0503 2509
119.0454 2981
123.0441 1145
136.0722 15907
147.0421 383
165.0539 225
216.0298 1549
```
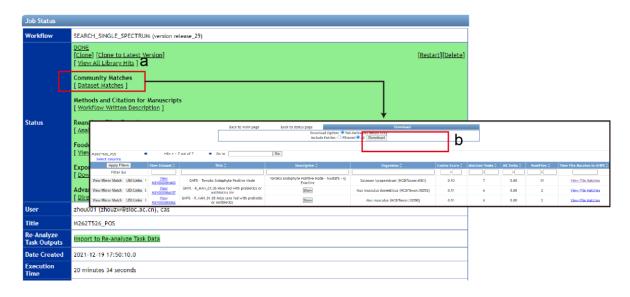
Modify the search parameters and click "submit" button. The **used parameters** in KGMN manuscript have been provided below.
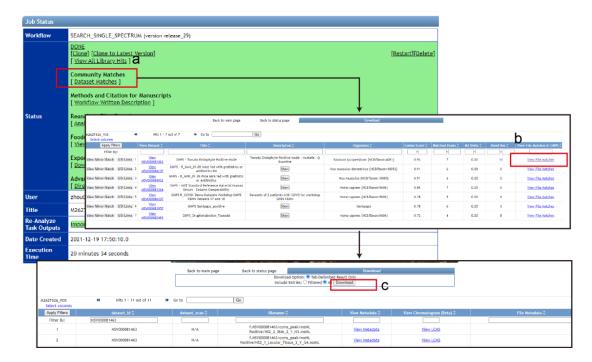


When the job finished, you will receive an email with a link. You can view and download results in the webserver.

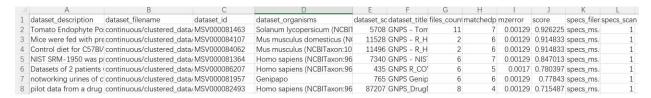- Matched data set: Dataset Matches → View File Matches → Download

- Matched files: Dataset Matches → View File Matches → Download



## 3. Result interpretation and visualization.

The downloaded results include 2 ZIP files, "view_all_datasets_matched.zip" and "view_all_file_datasets_matched.zip". The files in packages can be further opened with Microsoft Office Excel or other program tools (e.g. R, Python).

- The table of "view_all_datasets_matched" contains meta information of appeared data sets, like "dataset description", "dataset id", "dataset organisms" and "files count". Furthermore, we can conclude the species and sample information based on the dataset description. For our examples, it was appeared in 7 datasets, and 3 organisms (where genipapo is from human urine actually according to the data set description).

- The table of "view_all_file_datasets_matched" contains names of matched files. Each file can be viewed online through the filename in GNPS dashboard ([https://gnps-lcms.ucsd.edu/](https://gnps-lcms.ucsd.edu/)), while the files and dataset can be accessed in GNPS datasets ([https://gnps.ucsd.edu/ProteoSAFe/datasets.jsp](https://gnps.ucsd.edu/ProteoSAFe/datasets.jsp)).

| | A | B | C | D | E |
|---|---|---|---|---|---|
| 1 | basefilename | cluster_sca | dataset_id | filename | metadata |
| 2 | 018c.mzML | 435 | MSV000086207 | f.MSV000086207/ccms_peak/018c.mzML | |
| 3 | 018b.mzML | 435 | MSV000086207 | f.MSV000086207/ccms_peak/018b.mzML | |
| 4 | 018a.mzML | 435 | MSV000086207 | f.MSV000086207/ccms_peak/018a.mzML | |
| 5 | 017c.mzML | 435 | MSV000086207 | f.MSV000086207/ccms_peak/017c.mzML | |
| 6 | 017b.mzML | 435 | MSV000086207 | f.MSV000086207/ccms_peak/017b.mzML | |
| 7 | 017a.mzML | 435 | MSV000086207 | f.MSV000086207/ccms_peak/017a.mzML | |
| 8 | E12_3.mzML | 11528 | MSV000084107 | f.MSV000084107/ccms_peak/E12_3.mzML | |
| 9 | E12_2.mzML | 11528 | MSV000084107 | f.MSV000084107/ccms_peak/E12_2.mzML | |
| 10 | E12_3.mzML | 11496 | MSV000084062 | f.MSV000084062/ccms_peak/E12_3.mzML | |
| 11 | E12_2.mzML | 11496 | MSV000084062 | f.MSV000084062/ccms_peak/E12_2.mzML | |
| 12 | DM000088099_RB7_01_29 | 87234 | MSV000082493 | f.MSV000082493/ccms_peak/urine/DM000088099_RB | |
| 13 | DM000086580_RF12_01_2 | 87207 | MSV000082493 | f.MSV000082493/ccms_peak/urine/DM000086580_RF1 | |
| 14 | DM000078719_RA11_01_2 | 87214 | MSV000082493 | f.MSV000082493/ccms_peak/urine/DM000078719_RA | |
| 15 | DM000078708_RC10_01_2 | 87214 | MSV000082493 | f.MSV000082493/ccms_peak/urine/DM000078708_RC | |
| 16 | DM000078265_RD7_01_29 | 87207 | MSV000082493 | f.MSV000082493/ccms_peak/urine/DM000078265_RD | |
| 17 | DM000076834_RB8_01_29 | 87230 | MSV000082493 | f.MSV000082493/ccms_peak/urine/DM000076834_RB | |
| 18 | DM000076821_RC12_01_2 | 87234 | MSV000082493 | f.MSV000082493/ccms_peak/urine/DM000076821_RC | |
| 19 | DM000076799_RC8_01_29 | 87230 | MSV000082493 | f.MSV000082493/ccms_peak/urine/DM000076799_RC | |
| 20 | Urine83_Juice_12h_Top3_F | 765 | MSV000081957 | f.MSV000081957/ccms_peak/Urine83_Juice_12h_Top3_ | |

With above information, it would be easy to reproduce figures of repository validation. The result of above example can be downloaded **here**.

**Reference:**

1. Wang, M., Jarmusch, A.K., Vargas, F. et al. Mass spectrometry searches using MASST. Nat Biotechnol 38, 23–26 (2020). https://doi.org/10.1038/s41587-019-0375-9