
Qiita-GNPS-workshop

Release 0.01

May 08, 2017

1	Qiita tutorials:	3
1.1	Getting CMI Workshop example data	3
1.2	Setting up Qiita	3
1.3	Creating a test study	5
1.4	Creating an example study	5
1.5	Adding sample information	8
1.6	16S Data Processing in Qiita	11
1.7	16S Microbiome Analysis in Qiita	23
1.8	Notes on metabolomics	33
1.9	Feature finding with MZmine2	35
1.10	Metabolomics demo data in Qiita	46
1.11	GNPS tutorial for MS/MS data annotation	46

Materials below are intended for CMI Qiita/GNPS workshop participants. They include all information covered during days 1 and 2 of the workshop.

For more information on Qiita, including Qiita philosophy and documentation, please visit [Qiita website](#).

For general information about workshops, please [contact Tomasz Kosciolk](#) directly.

CHAPTER 1

Qiita tutorials:

This tutorial will walk you through creation of your account and a test study in Qiita.

Getting CMI Workshop example data

First, we'll download some [example data](#). These files contain both 16S and shotgun metagenomics data for 12 samples from the American Gut Project.

For this tutorial, the relevant files are:

```
qiita-files/16S/*.fastq.gz      # The actual 16S sequences, one per sample
qiita-files/sample_information.txt # The sample information file
qiita-files/prep_template_16S.txt # The prep information file
```

Next, we'll sign up for Qiita and create a study for these data.

Setting up Qiita

Signing up for a Qiita account

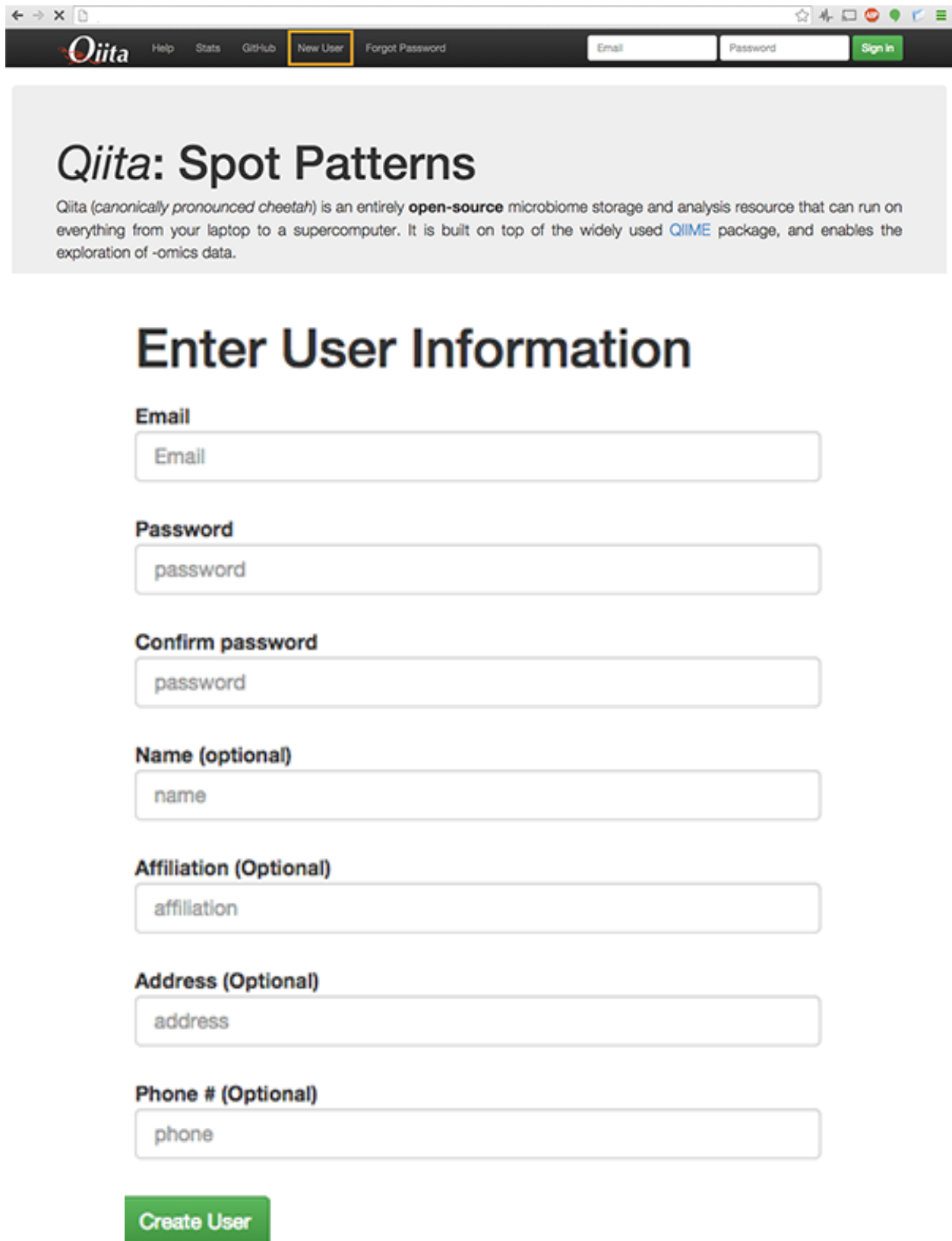
Open your browser (it must be Chrome or Firefox) and go to [Qiita \(https://qiita.ucsd.edu\)](https://qiita.ucsd.edu).

Click on “New User”.

The “New User” link brings you to a page on which you can create a new account. Optional fields are indicated explicitly, while all other fields are required. Once the form is submitted, an email will be sent to you containing instructions on how to verify your email address.

Logging into your account and resetting a forgotten password

Once you have created your account, you can log into the system by entering your email and password.



The screenshot shows a web browser window with the Qiita website. The navigation bar at the top includes the Qiita logo, links for Help, Stats, GitHub, New User (highlighted with a yellow box), and Forgot Password. There are also input fields for Email and Password, and a Sign In button. The main content area has a heading 'Qiita: Spot Patterns' followed by a paragraph describing Qiita as an open-source microbiome storage and analysis resource. Below this is a large heading 'Enter User Information' and a registration form with the following fields: Email, Password, Confirm password, Name (optional), Affiliation (Optional), Address (Optional), and Phone # (Optional). Each field has a placeholder text. At the bottom of the form is a green 'Create User' button.

Qiita: Spot Patterns

Qiita (*canonically pronounced cheetah*) is an entirely **open-source** microbiome storage and analysis resource that can run on everything from your laptop to a supercomputer. It is built on top of the widely used [QIIME](#) package, and enables the exploration of -omics data.

Enter User Information

Email

Password

Confirm password

Name (optional)

Affiliation (Optional)

Address (Optional)

Phone # (Optional)

Create User



If you forget your password, you will need to reset it. Click on “Forgot Password”.



This will take you to a page on which to enter your email address; once you click the “Reset Password” button, the system will send you further instructions on how to reset your lost password.

Lost Password

Enter email for account

Email

[Reset Password](#)

Updating your settings and changing your password

If you need to reset your password or change any general information in your account, click on your email at the top right corner of the menu bar to access the page on which you can perform these tasks.

Creating a test study

Studies are the source of data for Qiita. Studies can contain only one set of samples but can contain multiple sets of raw data, each of which can have a different preparation – for example, 16S, shotgun metagenomics, and metabolomics, or even multiple preparations of the same type (e.g., a plate rerun, biological and technical replicates, etc).

In this tutorial, our study contains 12 samples, each with two types of data: 16S and shotgun metagenomics. To represent this project in Qiita, you will need to create a single study with a single sample information file that contains all 12 samples. Then, you will link separate preparation files for each data type.

Creating an example study

To create a study, click on the “Study” menu and then on “Create Study”. This will take you to a new page that will gather some basic information to create your study.

The “Study Title” has to be unique system-wide. Qiita will check this when you try to create the study, and may ask you to alter the study name if the one you provide is already in use.

A principal investigator is required, and a list of known PIs is provided. If you cannot find the name you are looking for in this list, you can choose to add a new one.

Welcome antgonza@gmail.com

Log Out

User Information

Name

Affiliation

Address

Phone

Save Edits

Change Password

Old Password

New Password

Repeat New Password

Change Password

Study ▾

Help

Stats

Create Study

View My Studies

View Public Studies

Create a new Study

* = Required Field

Study Title	<input type="text" value="JGS April 2017 CMI workshop"/>	*
Study Alias	<input type="text" value="April 2017 CMI"/>	*
DOI	<input type="text"/>	
Just values, no links, comma separated values		
PUBMED ID	<input type="text"/>	
Just values, no links, comma separated values		
Study Abstract	<input type="text" value="An example study"/>	*
Study Description	<input type="text" value="Just a test"/>	*
Principal Investigator	<input type="text" value="Rob Knight, UCSD"/>	*
Lab Person	<input type="text" value="Jon Sanders, UCSD"/>	
	Can't find the person you're looking for? Add a person	
	<div> air built environment host-associated human-amniotic-fluid human-associated human-blood human-gut human-oral human-skin human-urine human-vaginal microbial mat/biofilm miscellaneous natural or artificial environment plant-associated sediment soil wastewater/sludge water </div>	*
Environmental Packages		
You can select multiple entries by control-clicking (mac: command-clicking)		
Event-Based Data	<input type="text" value="No timeseries"/>	
<input type="button" value="Create Study"/>		

Select the environmental package appropriate to your study. Different packages will request different specific information about your samples. This information is optional; for more details, see the metadata section.

There is also an option to specify time series type (“Event-Based Data”) if you have time series data. In our case, the samples come from a cross-sectional study design, so you should select “No time series.” For more information on time series types, you can check out the [in-depth tutorial](#) on the Qiita website.

Once your study has been created, you will be informed by a green message; click on the study name to begin adding your data.

The system has been updated. Please report any issue.
Study JGS April 2017 CMI workshop successfully created

Adding sample information

Sample information is the set of metadata that pertains to your biological samples: these are the measured variables that are motivating you to look for response variables in the microbiome. **IMPORTANT:** your metadata are your study; it is imperative that those data are consistent, correct, and sufficiently detailed. (To learn more, including how to format your own sample info file, check out the [in-depth documentation](#) on the Qiita website.)

The first point of entrance to a study is the study description page. Here you will be able to edit the study info, upload files, and manage all other aspects of your study.

The screenshot shows the Qiita web interface. At the top is a navigation bar with the Qiita logo and links for Analysis, Study, More Info, Current and Future Features, Downloads, and a user profile for jonsan@gmail.com. On the left sidebar, there are buttons for Study Information, Sample Information, Upload Files, and download links for QIIME maps and BIOMs, and all raw data. The main content area displays the study title 'JGS April 2017 CMI workshop - ID 10965' and a sub-header 'April 2017 CMI'. Below this is an 'Abstract' section with the text 'An example study'. Further down, study details are listed: Study ID: 10965, PI: Rob Knight (UCSD), Lab Contact: Jon Sanders (UCSD), Shared With, Samples: 0, and EBI: not submitted. There are 'Share', 'Edit', and 'Delete' buttons. To the right, there is a 'Study Tags' section with a list of tags: 'Previously admin', 'Previously assigned', and 'New'. Below the tags is an 'Add more tags' input field and a 'Save tags' button. A message at the bottom states 'No preparation information has been added yet'. At the very bottom, there is a footer with a thank you message and contact information.

The first step after study creation is uploading files. Click on the “Upload Files” button: as shown in the figure below, you can now drag-and-drop files into the grey area or simply click on “select from your computer” to select the fastq, fastq.gz or txt files you want to upload.

Uploads can be paused at any time and restarted again, as long as you do not refresh or navigate away from the page, or log out of the system from another page.

Drag the file named “sample_information.txt” into the upload box. It should upload quickly and appear with a check-box next to it below.

Uploading files for: JGS April 2017 CMI workshop (April 2017 CMI)

Currently we can process (fastq, fastq.gz, txt, tsv, sff, fasta, fna, qual, biom):

- fastq or fastq.gz (gzipped fastq) for sequences. Note that zip files can not be processed
- tab separated text files for sample and prep templates, the extension should be txt

[Go to study description](#)

Upload files (max file size: 2.0 TB)

Drop files here to upload or [select from your computer](#)

100 % ▶ ||

Keep track of your upload or pause/resume it! ↑
(you can even close your computer or change networks!)

sample_information.txt ☐

 Delete selected files

Once your file has uploaded, click on “Go to study description” and, once there, click on the “Sample Information” tab. Select your sample information from the dropdown menu next to “Upload information” and click “Create”.

JGS April 2017 CMI workshop - ID 10965

April 2017 CMI

Sample Information

Upload information: sample_information.txt ▾ Create

If uploading a qiime mapping file, select data type: Choose a data type... ▾

If something is wrong with the sample information file, Qiita will let you know with a red banner at the top of the screen.

If the file processes successfully, you should be able to click on the “Sample Information” tab and see a list of the imported metadata fields.

You can also click on “Sample Summary” to check out the different metadata values. Select a metadata column to visualize in the dropdown menu and click “Add column.”

In this cohort, only three people were sensible enough to own a cat.

Next, we’ll add 16S data and do a preliminary analysis.

❗ The 'sample_name' column is missing from your template, this file cannot be parsed.
Need help? Send us an [email](#).

JGS April 2017 CMI workshop - ID 10965

April 2017 CMI

Sample Information

Download Sample Info

Delete

There are **23** samples and **255** columns in this study.

Upload information:

Choose file...

Update

Sample information summary



vioscreen_frtsumm

values



vioscreen_fibinso

values



vioscreen_fried_fish_servings

values

JGS April 2017 CMI workshop - ID 10965

April 2017 CMI

Sample Summary

Add sample column information to table

cat

Add column

	Sample	cat
	10965.29511	No
	10965.31151	No
	10965.27689	No

Next: *16S Data Processing in Qiita*

16S Data Processing in Qiita

Now, we'll upload some actual microbiome data to explore. To do this, we need to add the data themselves, along with some information telling Qiita about how those data were generated.

Adding a preparation template and linking it to raw data

Where the *sample info file* has the biological metadata associated with your samples, the *preparation info file* contains information about the specific technical steps taken to go from sample to data. Just as you might use multiple data-generation methods to get data from a single sample – for example, target gene sequencing and shotgun metagenomics – you can have multiple prep info files in a single study, associating your samples with each of these data types. You can learn more about prep info files at the [Qiita documentation](#).

Go back to the “Upload Files” interface. In the [example data](#), find and upload the files in the *16S* folder and the file called *prep_information_16S.txt*.

Now you can click the “Add New Preparation” button. This will bring up the following dialogue:

JGS April 2017 CMI workshop - ID 10965

April 2017 CMI

Select file: *

Select data type: *

Select Investigation Type:

* Required fields

Create New Preparation

Select *prep_information_16S.txt* from the “Select file” dropdown, and *16S* as the data type. Optionally, you can also select one of a number of investigation types that can be used to associate your data with other like studies in the database. Click “Create New Preparation”.

You should now see a summary of your preparation info, similar to the summary we saw of the sample info:

In addition, you should see a “16S” button appear under “Data Types” on the menu to left:

You can click this to reveal the individual prep info files of that data type that have been associated with this study:


If you have multiple 16S preparations (for example, if you sequenced using several different primer sets), these would each show up as a separate entry here.

Now, you can associate the sequence data from your study with this preparation.

JGS April 2017 CMI workshop - ID 10965

April 2017 CMI

Prep information 2712 - ID 2712 (16S)

 Prep info

 QIIME map

 Delete

There are **23** samples and **36** columns in this preparation.

Update information:


Select Investigation Type:


No files attached to this preparation

Select type:

Add a name for the file:

Information summary

 **barcode:** All the values in this category are different.

 **center_name:** UCSDMI is repeated in all rows.

 **center_project_name**


AG21

12

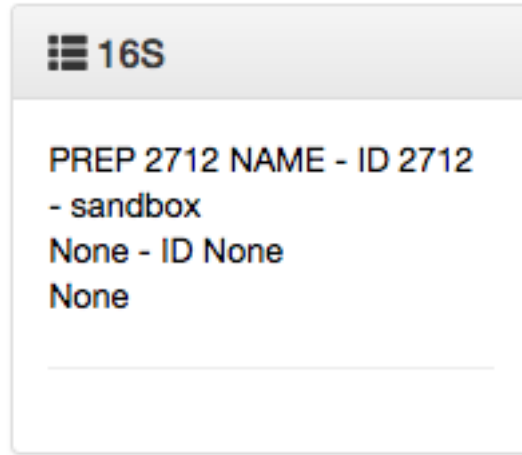
AG22

11

Data Types (click on the tabs)

 16S

Data Types (click on the tabs)



In the prep info dialogue, there is a dropdown menu below the words *No files attached to this preparation*, labeled “Select type”. Click “Choose a type” to see a list of available file types. In our case, we’ve uploaded FASTQ-formatted files per each sample in our study, so we will choose *per_sample_FASTQ*.

Magically, this will prompt Qiita to associate your uploaded files with the corresponding samples in your preparation info. (Our prep info file has a column named *run_prefix*, which associated the *sample_name* with the file name prefix for that particular sample.)

You should see this as a list of filenames showing up in the green *raw forward seqs* column below the import dropdown. You’ll want to give the set of these per-sample-FASTQ files a name (**Add a name for the file**), and then click “Add files” below.

That’s it! Your data are ready for processing.

Exploring the raw data

Click back through on your 16S preparation. Now that you’ve associated sequence files with this prep, you’ll have a *Files network* displayed:

Your collection of FASTQ files for this prep are all represented by a single object in this network, currently called *dflt_name*. Click on the object.

Now, you’ll have a series of choices for interacting with this object. You can click “Edit” to rename the object, “Process” to perform analyses, or “Delete” to delete it. In addition, you’ll see a list of the actual files associated with this object.

Scroll to the bottom, and you’ll also see an option to generate a summary of the object.

If you click this button, it will be replaced with a notification that the summary generation has been added to the processing queue.

To check on the status of the processing job, you can click the rightmost icon at the top of the screen:

There are **23** samples and **36** columns in this preparation.

Update information:

Select Investigation Type:

No files attached to this preparation

Select type: ☒ Choose a type...

Add a name

- BIOM - BIOM table
- Demultiplexed - Demultiplexed and QC sequences
- FASTA - None
- FASTA_Sanger - None
- FASTQ - None
- per_sample_FASTQ - None**
- SFF - None

Information



No files attached to this preparation

Select type:

Add a name for the file:

Now, you can import files from other studies

or click and drag your uploaded files to the correct file type

Please make sure that the correct files are in the correct column.

Note: the system will try to auto select the files based on run_prefix, if that doesn't work, either the type you selected doesn't support the use of run_prefix or the run_prefix is wrong

Available Files

raw forward seqs

raw reverse seqs

F000010242.small.1

F000018425.small.1

F000030023.small.1

F000018424.small.1

F000030230.small.1

Add files

Information summary



barcode: All the values in this category are different.

Files network

(Click nodes for more information, blue are jobs)

Check our [data processing recommendations](#).



This will open a dialogue that gives you information about currently running jobs, as well as jobs that failed with some sort of error.

The summary generation shouldn't take too long. When it completes, you can click back on the `per_sample_FASTQ` object and scroll to the bottom of the page to see a short peek at the data in each of the FASTQ files in the object. These summaries can be useful for troubleshooting.

Now, we'll process the raw data into something more interesting.

Processing 16S data

Scroll back up and click on the `per_sample_FASTQ` object, and select "Process". This will bring you to another network visualization interface. Here, you can add processing steps to your objects.

Click again on the `per_sample_FASTQ` object. Below the files network, you will see an option to *Choose command*. Based on the type of object, this dropdown menu will give you a list of available processing steps.

For 16S `per_sample_FASTQ` objects, the only available command is *Split libraries FASTQ*. This converts the raw FASTQ data into the file format used by Qiita for further analysis (you can read more extensively about this file type [here](#)).

Select the *Split libraries FASTQ* step. Now, you will be able to select the specific combination of parameters to use for this step in the *Choose parameter set* dropdown menu.

For our files, choose *per sample FASTQ defaults, phred_offset 33*. The specific parameter values used will be displayed below. (The other commonly used choice for data generated at the CMI is *golay_12, reverse complement mapping file barcodes, reverse complement barcodes*, which is what you will select if you have one set of non-demultiplexed FASTQ files (forward, reverse, and barcode) containing all of your samples.)

Click "Add Command".

You'll see the files network update. In addition to the original grey object, you should now see the processing command (represented in blue) and the object produced from that command (also represented in grey).

Files network


(Click nodes for more information, blue are jobs)

Check our data [processing recommendations](#).


dflt_name - per_sample_FASTQ

dflt_name (ID: 26007)

 Edit

 Process


 Delete


Processing parameters:


Visibility: sandbox


[Request approval](#)

Available files:

 F000010242.small.fastq.gz (raw forward seqs)

 F000018425.small.fastq.gz (raw forward seqs)

 F000030023.small.fastq.gz (raw forward seqs)

 F000030230.small.fastq.gz (raw forward seqs)

Currently, no summary exists.

[Generate summary](#)

Log Out

 (6)





Processing Jobs

(skipping successful jobs)

Search:

	Heartbeat ▼	Name ▲	Status ▲	Step
	2017-04-18 20:17:12	Generate HTML summary	running	
	2017-02-14 07:54:13	Validate	error	Step 2: Validating 'per_sample_FASTQ' files

↓ F000030230.small.fastq.gz (raw forward seqs)

↓ artifact_26010.html (html summary)

F000010242.small.fastq.gz (raw_forward_seqs)

MD5: 6370d7049892c401037c72161c52b0e8

@10317.000010242_58737 orig_bc=GCGTCTGTAGCA new_bc=ACGTCTGTAGCA

bc_diffs=1

TACGGAGGGTGCAAGCGTTAATCGGAATTACTGGGCGTAAAGCGCACGCAGGCGGTTTGTGAAGTCAGATGTGAAATC

+

BBBBBB@10317.000010242_70616 orig_bc=TCGTCTGTAGCA new_bc=ACGTCTGTAGCA

bc_diffs=1

TACGTAGGGGGCGAGCGTTATCCGGATTTCATTGGGCGTAAAGCGCGCGTAGGCGGCCCCGGCAGGCCGGGGGTCGAAGC

+

BBBB@10317.000010242_2480 orig_bc=ACGTCTGTAGCA new_bc=ACGTCTGTAGCA

bc_diffs=0

TACAGAGGGTGCAAGCGTTAATCGGAATTACTGGGCGTAAAGCGCGCGTAGGTGGTTTGTGAAGTTGAATGTGAAATC

F000018425.small.fastq.gz (raw_forward_seqs)

MD5: a7e9451e53b580b17eb9204c0731acd7

@10317.000018425_61392 orig_bc=TTATGGGTAAAGA new_bc=TTATGGGTAAAGA

bc_diffs=0

Processing dflt_name (ID: 26010)

Processing workflow

▶ Run

Don't forget to hit "Run" once you are done with your workflow!

Wondering what to select? Check our data [processing recommendations](#).



dflt_name - (per_sample_FASTQ)

Choose command: Choose command... ▾

Choose command: Split libraries FASTQ ▾

input data: dflt_name ▾

Choose parameter set: per sample FASTQ defaults

max_barcode_errors 1.5

barcode_type not-barcoded

max_bad_run_length 3

phred_offset auto

rev_comp false

phred_quality_threshold 3

rev_comp_barcode false

rev_comp_mapping_barcode false

min_per_read_length_fraction 0.75

sequence_max_n 0

Add Command



You can click on the command to see the parameters used, or on an object to perform additional steps.

Note that the command hasn't actually been run yet! (We'll still need to click "Run" at the top.) This allows us to add multiple processing steps to our study and then run them all together.

We're going to process our sequences files using two different workflows. In the first, we'll use a conventional reference-based OTU picking strategy to cluster our 16S sequences into OTUs. This approach matches each sequence to a reference database, ignoring sequences that don't match the reference. In the second, we will use [deblur](#), which uses an algorithm to remove sequence error, allowing us to work with unique sequences instead of clustering into OTUs. Both of these approaches work great with Qiita, because we can compare the observations between studies without having to do any sort of re-clustering!

The closed reference workflow

To do closed reference OTU picking, click on the *demultiplexed* object and select the *Pick closed-reference OTUs* command. We will use the *default - serial* parameter set for our data, which are relatively small. For a larger data set, we might want to use the parallel implementation.

By default, Qiita uses the GreenGenes 16S reference database. You can also choose to use Silva, or the Unite fungal ITS database.

Click "Add Command", and you will see the network update:



Here you can see the blue "Pick closed-reference OTUs" command added, and that the product of the command is a BIOM-formatted OTU table.

That's it!

The deblur workflow

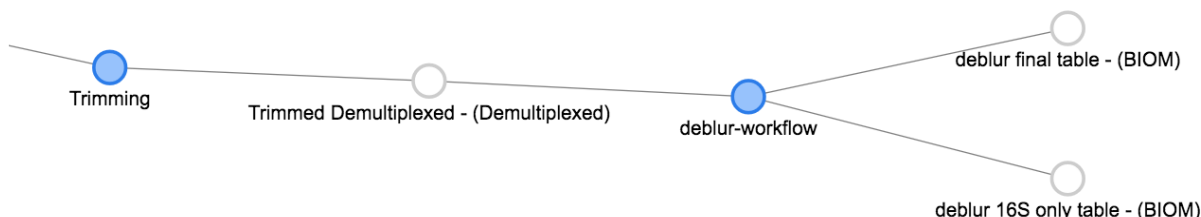
The deblur workflow is only marginally more complex. Although you can deblur the demultiplexed sequences directly, *deblur* works best when all the sequences are the same length. By trimming to a particular length, we can also ensure our samples will be comparable to other samples already in the database.

Click back on the *demultiplexed* object. this time, select the *Trimming* operation. Currently, there are three trimming length options. Let's choose *Trimming 100*, which trims to the first 100bp, for this run, and click "Add Command".

Now you can see that we have the same *demultiplexed* object being used for two separate processing steps – closed-reference OTU picking, and trimming.



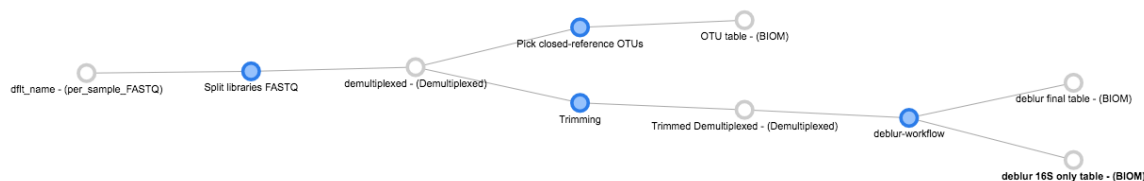
Now we can click the *Trimmed Demultiplexed* object and add a *deblur* step. Choose *deblur-workflow* from the *Choose command* dropdown, and *Defaults* for the parameter set. Add this command.



As you can see, *deblur* produces two BIOM-formatted OTU tables as output. The *deblur 16S only table* contains deblurred sequences that have been filtered to try and exclude things like organellar mitochondrial reads, while *deblur final table* has all the sequences.

Running the workflow

Now, we can see the whole set of commands and their output files:



Click “Run” at the top of the screen, and Qiita will start executing all of these jobs. You’ll see a “Workflow submitted” banner at the top of your window.

As noted above, you can follow the process of your commands in the dialogue at the top right of the window.

You can also click on the objects in the prep info file network, and see status updates from the commands running on that object at the bottom of the page:

Once objects have been generated, you can generate summaries for them just as you did for the original *per_sample_FASTQ* object.

The summary for the *demultiplexed* object gives you information about the length of sequences in the object:

The summary for a BIOM-format OTU table gives you a histogram of the the number of sequences per sample:

Next: *16S Microbiome Analysis in Qiita*

Jobs using this set of files:

Job **2915b8ad-b5a3-4e55-8179-d07f0b61595f** (Pick closed-reference OTUs).

Status: *running*. Step: *Step 3 of 4: Executing OTU picking*

Currently, no summary exists.

Features

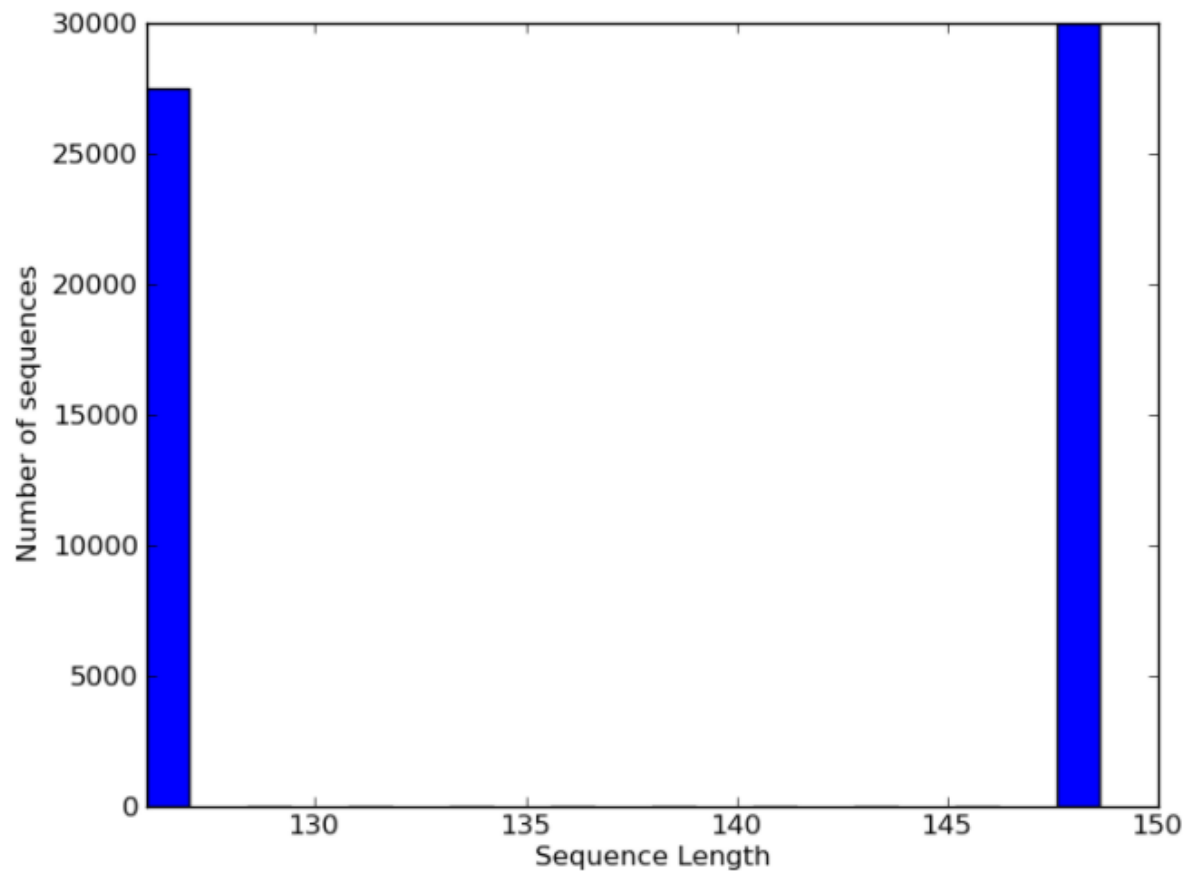
Total: 57500

Max: 150

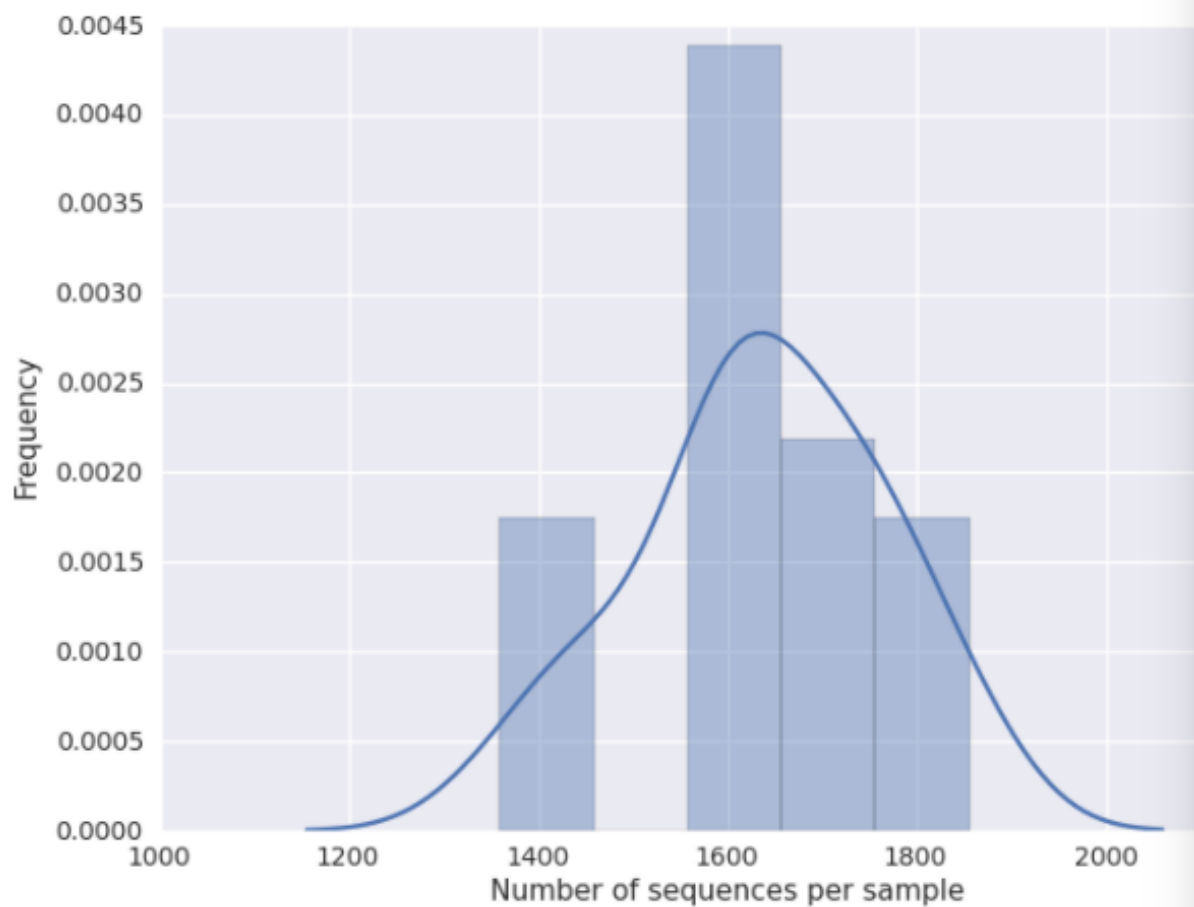
Mean: 138

Standard deviation: 150

Median: 11



Number of samples: 23
Number of features: 207
Minimum count: 1359
Maximum count: 1855
Median count: 1631
Mean count: 1635




16S Microbiome Analysis in Qiita

Analysis of Closed Reference processing

To create an analysis, select *Create new analysis* from the top menu.

This will take you to a list of studies with samples available to you for analysis, divided between your studies and publically available (‘Other’) studies.


Analysis - Study - More Info - Current and Future Features Downloads -
Welcome jonsan@gmail.com Log Out

Filter studies by tags: (Admin, User)

GOLD NASH HFD water polar hypersaline Pregnancy TLR3 Mouse Autism 16S hellbender mhc alleganiensis bishopi cancer
 MultipleSclerosis cryptobranchus pregnancy microbiota 16s pseudopregnancy gut

Your Studies (includes shared with you)

Show 5 entries

Filter results by column data (Title, abstract, PI, etc):

Expand	Add to analysis	Title	Study ID	Samples	Shared With These Users	Principal Investigator	Publications	Status	EBI
▼	No Processed Data	Baleen whales host a unique gut microbiome with similarities to both carnivores and herbivores	10285	0	Modify	Peter Girguis		sandbox	not submitted
▼	Add to Analysis	Lab contamination_Vibrio fischeri LOD plate extraction	10354	96	amnonim@gmail.com, Demo, Gail Ackermann, Jon Sanders Modify	Rob Knight		private	not submitted
▼	Add to Analysis	CAICE_Michaud	10356	43	amnonim@gmail.com, Jake, Jeff Mindrebo, Jon Sanders, zech Modify	Jennifer Michaud		private	not submitted
▼	Add to Analysis	Well-to-Well contamination	10401	768	amnonim@gmail.com, Jake, Jon Sanders Modify	Rob Knight		private	not submitted
▼	Add to Analysis	Knight_lab_master_mix_comparison_2016	10708	576	amnonim@gmail.com, James, Jon Sanders Modify	Rob Knight		private	not submitted

Showing 1 to 5 of 13 entries

Previous 1 2 3 Next

Other Studies

Expand	Add to analysis	Title	Study ID	Samples	Principal Investigator	Publications	EBI
▼	Add to Analysis	A core gut microbiome in obese and lean twins.	77	281	Jeff Gordon	19043404, 19043404	not submitted
▼	Add to Analysis	Soil bacterial and fungal communities across a pH gradient in an arable soil	94	27	Noah Fierer	20445636, 20445636	not submitted
▼	Add to Analysis	Succession of microbial consortia in the developing infant gut microbiome	101	63	Ruth Ley	20668239, 20668239	not submitted

Find the study you created for this tutorial under “Your Studies”. Click the down arrow at the left of the row. This will expand the study to expose all the objects from that study that are available to you for analysis.

▼	Add to Analysis	JGS April 2017 CMI workshop	10965	23	Modify	Rob Knight		sandbox	not submitted
---	---------------------------------	---	-------	----	------------------------	------------	--	---------	---------------

Processed Data

ID	Data type	Processed Date	Samples	Algorithm	Parameters
Add 26016	16S	2017-04-18 21:27:23.608563	23	deblur (deblur-workflow)	jobs-to-start: 5 pos-ref-db-fp: /databases/gg/13_8/sortmerna/88_otus threads-per-sample: 1 indel-prob: 0.01 neg-ref-fp: default trim-length: 100

You could add all of these objects to the analysis by selecting the *Add to Analysis* button. We will just add the Closed Reference OTU table object by clicking *Add* in that row.

Add	26019	16S	2017-04-18 22:39:24.764514	23	QIIME (Pick closed-reference OTUs)	<i>similarity:</i> 0.97 <i>reference_name:</i> Greengenes <i>sortmerna_e_value:</i> 1 <i>sortmerna_max_pos:</i> 10000 <i>threads:</i> 1 <i>sortmerna_coverage:</i> 0.97 <i>reference_version:</i> 13_8-97
-----	-------	-----	-------------------------------	----	------------------------------------	---

Now, the second-right-most icon at the top bar should be green, indicating that there are samples selected for analysis.



Clicking on the icon will take you to a page where you can refine the samples you want to include in your analysis. Here, all 23 of our samples are currently included:

Selected Samples

Create Analysis Clear Selected

JGS April 2017 CMI workshop

Processed Data

id	Datatype	Processed Date	Algorithm	Parameters	Samples		
26019	16S	2017-04-18 22:39:24.764514	QIIME (Pick closed-reference OTUs)	<i>similarity:</i> 0.97 <i>reference_name:</i> Greengenes <i>sortmerna_e_value:</i> 1 <i>sortmerna_max_pos:</i> 10000 <i>input_data:</i> 26014 <i>threads:</i> 1 <i>sortmerna_coverage:</i> 0.97 <i>reference_version:</i> 13_8-97	23	Show/Hide samples	Remove

You could optionally exclude particular samples from this set by clicking on “Show/Hide samples”, which will show each individual sample name along with a “remove” option. (Removing them here will mask them from the analysis, but will not affect the underlying files in any way.)

This should be good for now. Click the “Create Analysis” button, enter a name and description, then click “Create analysis”.

This brings you to the analysis commands selection page, where you can specify the steps in your analysis.

For this analysis, let’s go ahead and select the commands Summarize Taxa and Beta Diversity (Alpha Rarefaction can take some time to run).

We will also need to specify an even sampling or rarefaction depth. All the samples in the analysis will be randomly subsampled to this number of sequences, reducing potential biases. Samples with fewer than this number of sequences will be excluded, which can also be useful for excluding things like blanks.

You can get a good idea of where to set this threshold by looking at the histogram generated by summarizing the input closed-reference OTU table, as discussed in *16S Microbiome Analysis in Qiita*. Here, it looks like 2100 would be an appropriate cutoff: it excludes one clear outlier, but retains most of the samples.

Create new analysis ×

Analysis name

Description

Create analysis

Select Commands

Rarefaction Depth:

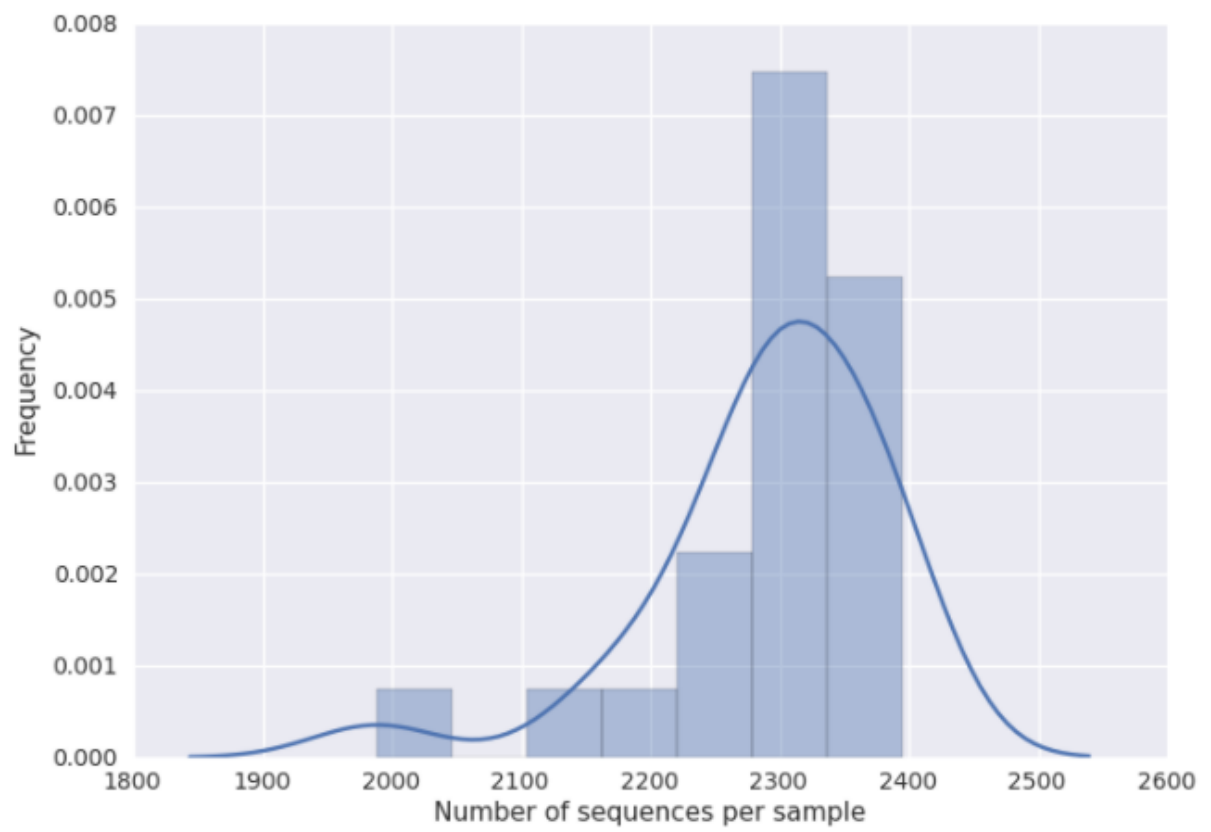
☐ Merge same sample ids

Merging sample ids is useful for when you have the same sample in different preparations of the same data, i.e. a sample processed twice in 16S. When the samples are not merged, they are prefixed with the artifact id.

Command

- ☒ Summarize Taxa
- ☒ Beta Diversity
- ☐ Alpha Rarefaction

Start Processing



Enter 2100 in the rarefaction depth field, select the check boxes for Summarize Taxa and Beta Diversity, and click “Start Processing”. You will see a list each step in the analysis, followed by its status:

Analysis Closed Ref Example

Thank you for using Qiita. For citations point to <http://qiita.microbio.me>.
Questions? qiita.help@gmail.com
Read our [terms and conditions](#).

Remove	16S: Beta Diversity: Running
Remove	16S: Summarize Taxa: Running
Remove	16S: Alpha Rarefaction: Running
Remove	tgz_analysis_10894: Queued
Remove	Finalize analysis: Queued
Remove	Creating Closed Ref Example... When finished, please click the 'Success' link to the right: Running

When the analysis is finished, click the ‘Success’ link to see the results.

The results page will have sections indication which samples were dropped due to insufficient numbers of reads, as well as sections for each data type.

Here, we have taxonomy summaries and beta diversity PCoA plots available.

Clicking on *bar_charts.html* under “Summarize Taxa” will take you to a visualization of the taxa that were found in your sample:

Under “Beta Diversity”, you will have a selection of Principle Coordinates Analyses of different measures of beta diversity, or the similarity between samples.

Clicking on one (say, *unweighted unifrac emperor pcoa plot*) will open an interactive visualization of the similarity among your samples. Generally speaking, the more similar the samples, the closer they are likely to be in the PCoA ordination. The Emperor visualization program offers a very useful way to explore how patterns of similarity in your data associate with different metadata categories. Here, I’ve colored the points in our test data by cat ownership.

Let’s take a few minutes now to explore the various features of Emperor. Open a new browser window with the [Emperor tutorial](#) and follow along with your test data.

Finally, if you ran Alpha Rarefaction, you will also have a link to interactive plots that can be used to show how different measures of alpha diversity correlate with different metadata categories:

Analysis of deblur processing

Creating an analysis of your deblurred data is virtually the same as the process for the Closed Reference data, but there are a few quirks.

First, because the deblur process creates two separate BIOM tables, you’ll want to make a note of the specific object ID number for the artifact you want to use. In my case, that’s ID 26017, the deblurred table with ‘only-16s’ reads.

The specific ID for your table will be unique, so make a note of it, and you can use it to select the correct table for analysis.

Analysis: CMI Workshop Test



Shared with:

Dropped Samples

16S

JGS April 2017 CMI workshop:

Total dropped: 1

10965.000030083

16S

Summarize Taxa

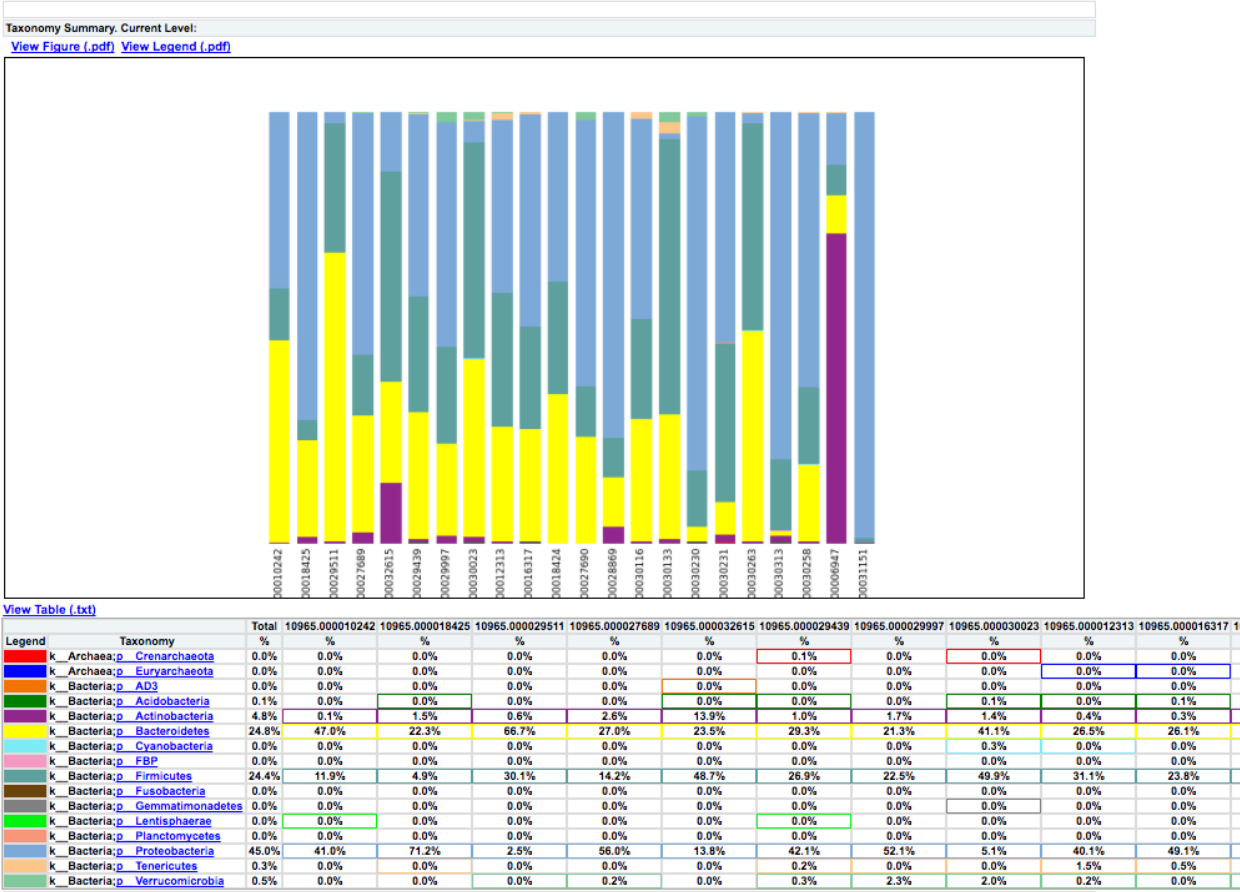
[area_charts.html](#)

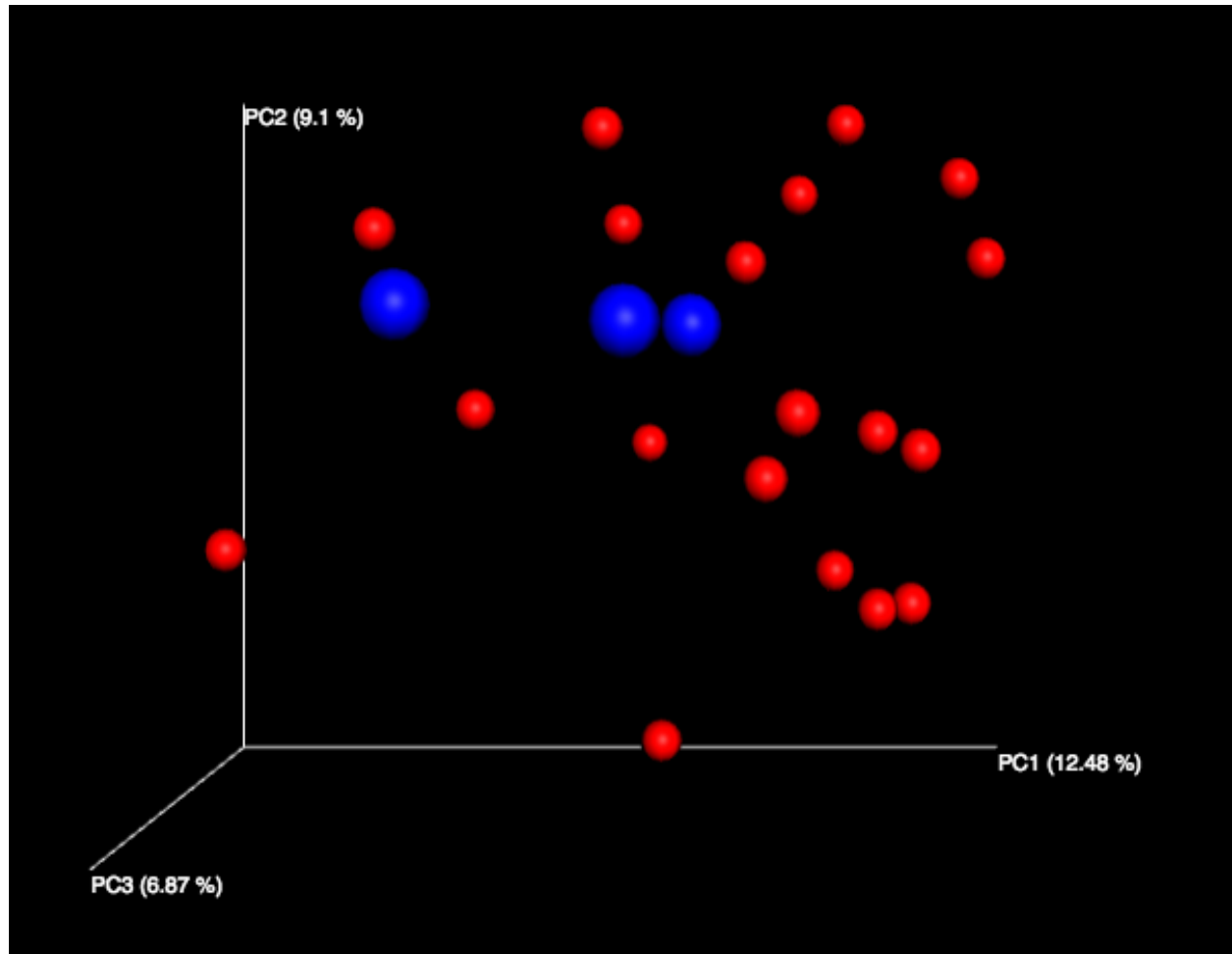
[bar_charts.html](#)

Beta Diversity

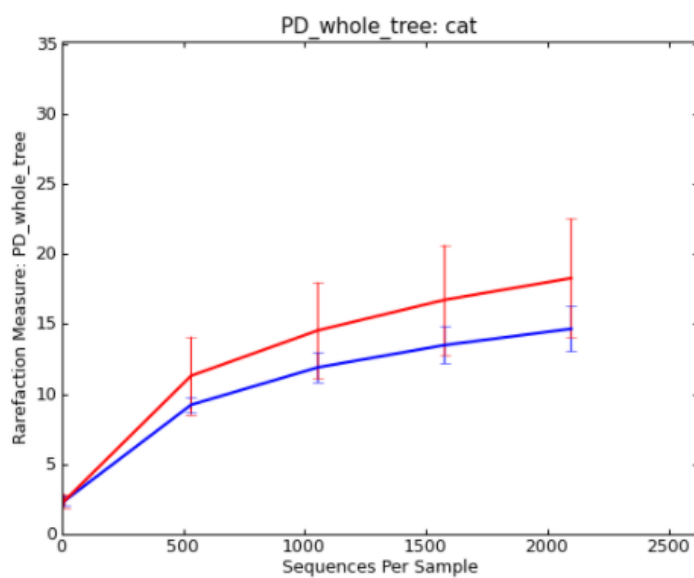
[unweighted unifracs emperor pcoa plot](#)

[weighted unifracs emperor pcoa plot](#)



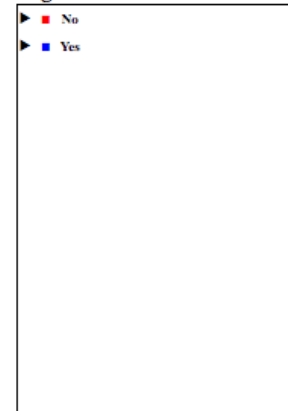


Select a Metric: Select a Category:



Show Categories:

Legend



dflt_name (ID: 26017)

Edit

Process

Delete


Processing parameters: jobs-to-start: 5 pos-ref-db-fp: /databases/gg/13_8/sortmerna/88_otus threads-per-sample: 1 indel-prob: 0.01 neg-ref-fp: default indel-max: 3 mean-error: 0.005 error-dist: 1, 0.06, 0.02, 0.02, 0.01, 0.005, 0.005, 0.005, 0.001, 0.001, 0.001, 0.0005 neg-ref-db-fp: default seqs-fp: 26015 skip-trimming: True negate: True pos-ref-fp: /databases/gg/13_8/rep_set/88_otus.fasta trim-length: 100 min-reads: 10 min-size: 2


Visibility: sandbox

Request approval

Available files:

 final.only-16s.biom (biom)

 final.seqs.fa.no_artifacts (preprocessed fasta)

 final.only-16s.biom.html (html summary)
Number of samples: 23**Number of features:** 207**Minimum count:** 1359**Maximum count:** 1855**Median count:** 1631**Mean count:** 1635

Second, currently only the Beta Diversity analysis command option is working with deblurred data.

Creating a meta-analysis

One of the most powerful aspects of Qiita is the ability to compare your data with hundreds of thousands of samples from across the planet. Right now, there are almost 130,000 samples publicly available for you to explore:

(You can get up-to-date statistics by clicking “Stats” under the “More Info” option on the top bar.)

Creating a meta-analysis is just like creating an analysis, except you choose data objects from multiple studies. Let’s start creating a meta-analysis by adding our Closed Reference OTU table to a new analysis.

Next, we’ll look for some additional data to compare against.

You noticed the ‘Other Studies’ table below ‘Your Studies’ when adding data to the analysis. (Sometimes this takes a while to load - give it a few minutes.) These are publicly available data for you to explore, and each should have processed data suitable for comparison to your own.

There are a couple tools provided to help you find useful public studies.

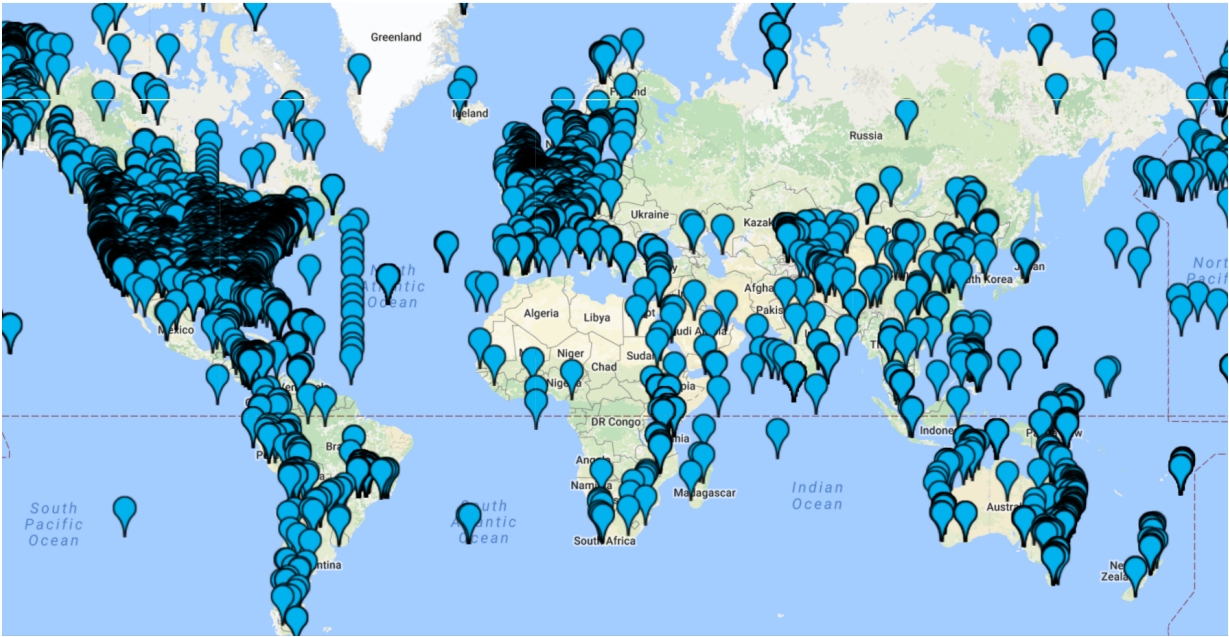
First, there are a series of “tags” listed at the top of the window:

There are two types of tags: admin-assigned (yellow), and user-assigned (blue). You can tag your own study with any tag you’d like, to help other users find your data. For some studies, Qiita administrators will apply specific reserved tags to help identify particularly relevant data. The “GOLD” tag, for example, identifies a small set of highly-curated, very well-explored studies. If you click on one of these tags, all studies not associated with that tag will disappear from the tables.

Second, there is a search field that allows you to filter studies in real time. Try typing in the name of a known PI, or a particular study organism – the thousands of publicly available studies will be filtered down to something that is easier to look through.

Generated on: 04-19-17 03:03:51

Studies	Samples	Users
<i>sandbox</i> : 793	<i>sandbox</i> : 130,583	
<i>public</i> : 310	<i>public</i> : 129,997	
<i>private</i> : 311	<i>private</i> : 115,951	2,978
<i>submitted to EBI</i> : 184	<i>submitted to EBI</i> : 107,888	
	<i>submitted to EBI (prep)</i> : 84,039	



Filter studies by tags: (Admin, User)

GOLD

NASH

HFD

water

polar

hypersaline

Pregnancy

TLR3

Mouse

Autism

16S

hellbender

mhc

alleganiensis

bishopi

cancer

MultipleSclerosis

cryptobranchus

pregnancy

microbiota

16s

pseudopregnancy

gut

Your Studies (includes shared with you)


Show 5 entries

Filter results by column data (Title, abstract, PI, etc):

Expand	Add to analysis	Title	Study ID	Samples	Shared With These Users	Principal Investigator	Publications	Status	EBI
--------	-----------------	-------	----------	---------	-------------------------	------------------------	--------------	--------	-----

Let's try comparing our data to the "Global Gut" dataset of human microbiomes from the US, Africa, and South America from the study "[Human gut microbiome viewed across age and geography](#)" by Yatsunenko et al. We can search for this dataset using the DOI from the paper: *10.1038/nature11053*.

Other Studies

Expand	Add to analysis	Title	Study ID	Samples	Principal Investigator	Publications	EBI
▼	Add to Analysis	 Human gut microbiome differentiation viewed across cultures, ages and families illumina	850	528	Jeff Gordon	22699611, 10.1038/nature11053	not submitted

Processed Data

ID	Data type	Processed Date	Samples	Algorithm	Parameters
Add 2458	16S	2015-10-15 15:14:56.647454	528	QIIME (Pick closed-reference OTUs)	similarity: 0.97 reference_name: Greengenes sortmerna_e_value: 1 sortmerna_max_pos: 10000 threads: 5 sortmerna_coverage: 0.97 reference_version: 13.8-97

Add the closed reference OTU table from this study to your analysis. You should now be able to click the green analysis icon in the upper right and see both your own OTU table and the public study OTU table in your analysis staging area:

You can now click "Create Analysis" just as before to begin specifying analysis steps. This time, let's just do the beta diversity step. Select the *Beta Diversity* command, enter a rarefaction depth of 2100, and click "Start Processing".

Because you've now expanded the number of samples in your analysis by more than an order of magnitude, this step will take a little longer to complete. But when it does, you will be able to use Emperor to explore the samples in your test dataset to samples from around the world!

Notes on metabolomics

Edited for the Dorrestein Lab by Louis-Felix Nothias, Daniel Petras and Ricardo Silva on December 2016. Last edit on April 2017.

About the metabolomics workshop

In the following documentation, we are providing step-by-step tutorials to perform basic analysis of liquid chromatography coupled to tandem mass spectrometry data (LC-MS/MS). These tutorials can be employed to process untargeted metabolomics data, such as those generated for seed funded project.

- The GNPS web-platform will be used to generate a qualitative analysis of your sample LC-MS/MS data. Such as the annotation of known compounds (by MS/MS spectral matching with public library), along as annotating unknown compounds by molecular networking (by spectral similarity).
- And we will use MZmine2 to process LC-MS/MS data in order to generate a feature table. This feature table contains the list of detected compounds and their relative distribution across samples. This feature table will be used to generate statistical analysis in Qiita.

Selected Samples

Create Analysis Clear Selected

JGS April 2017 CMI workshop

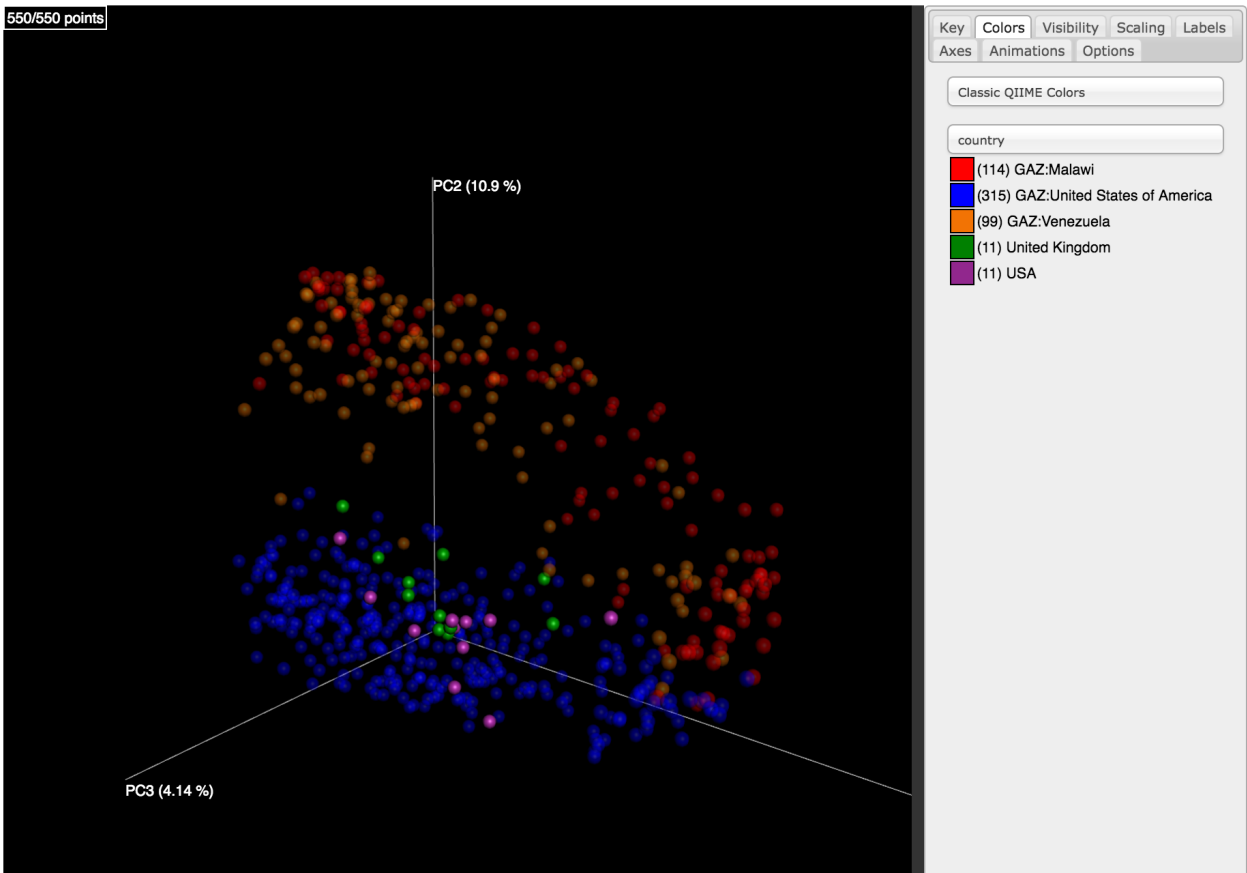
Processed Data

id	Datatype	Processed Date	Algorithm	Parameters	Samples		
26019	16S	2017-04-18 22:39:24.764514	QIIME (Pick closed-reference OTUs)	similarity: 0.97 reference_name: Greengenes sortmerg_e_value: 1 sortmerg_max_pos: 10000 input_data: 26014 threads: 1 sortmerg_coverage: 0.97 reference_version: 13_8-97	23	Show/Hide samples	Remove

Human gut microbiome differentiation viewed across cultures, ages and families
illumina

Processed Data

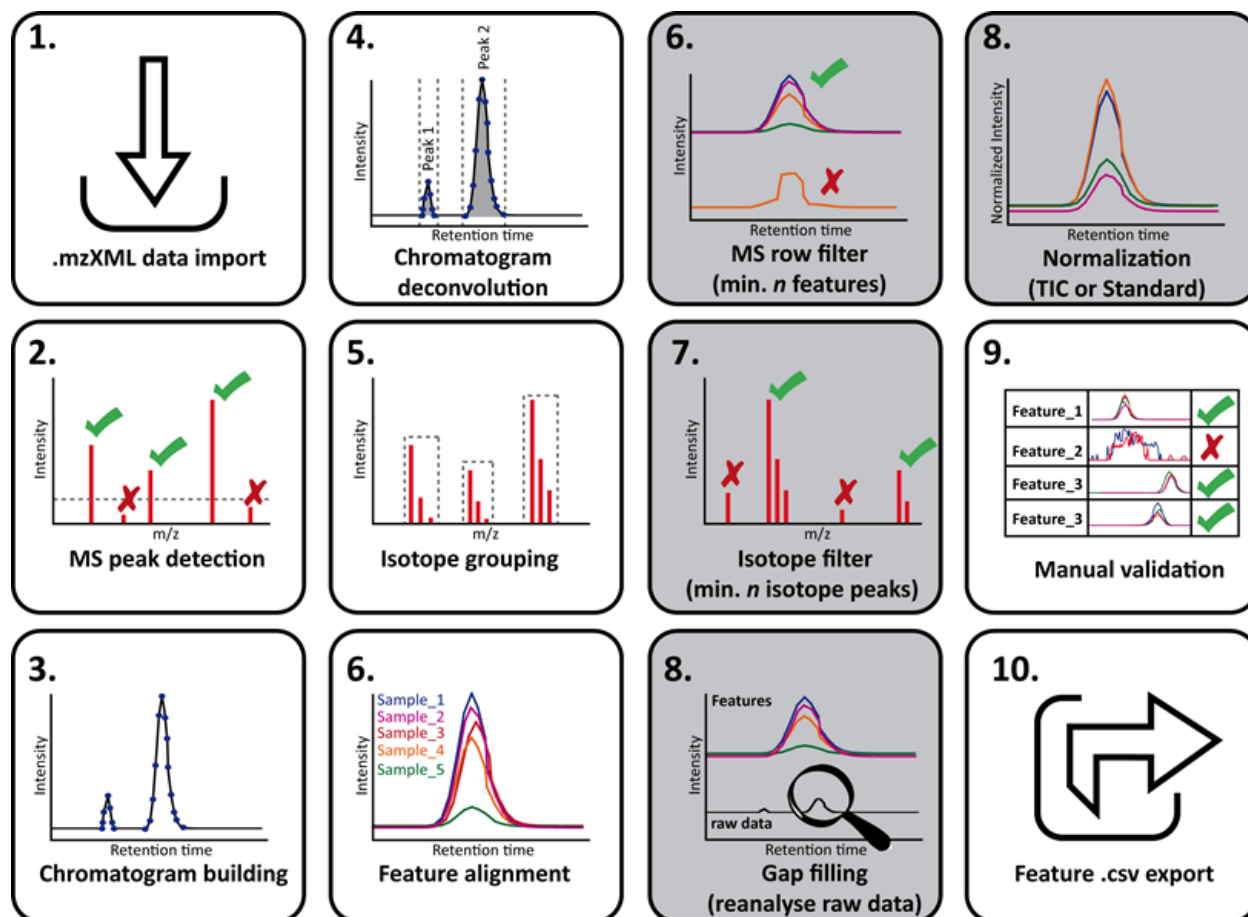
id	Datatype	Processed Date	Algorithm	Parameters	Samples		
2458	16S	2015-10-15 15:14:56.647454	QIIME (Pick closed-reference OTUs)	similarity: 0.97 reference_name: Greengenes sortmerg_e_value: 1 sortmerg_max_pos: 10000 input_data: 171 threads: 5 sortmerg_coverage: 0.97 reference_version: 13_8-97	528	Show/Hide samples	Remove



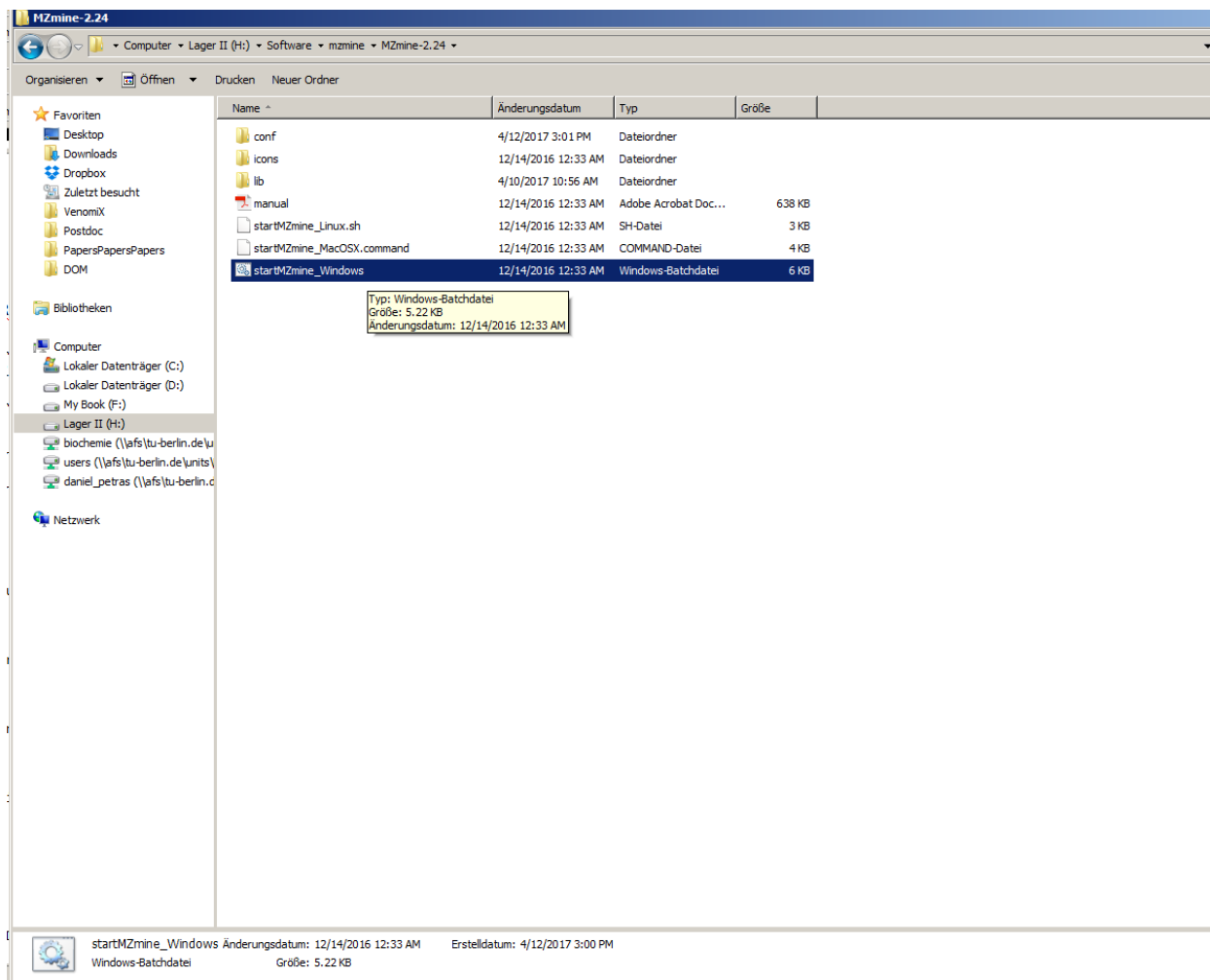
Feature finding with MZmine2

Please follow this [link](#) to install the software and dependencies.

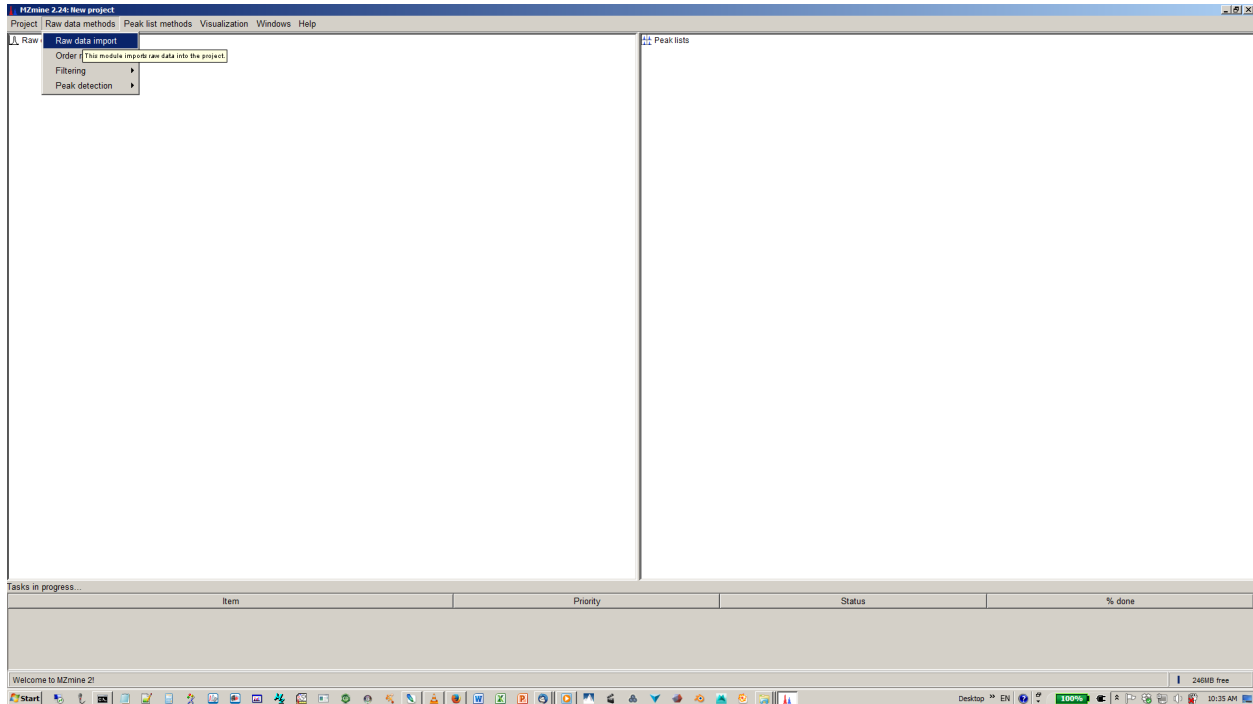
Complete workflow view



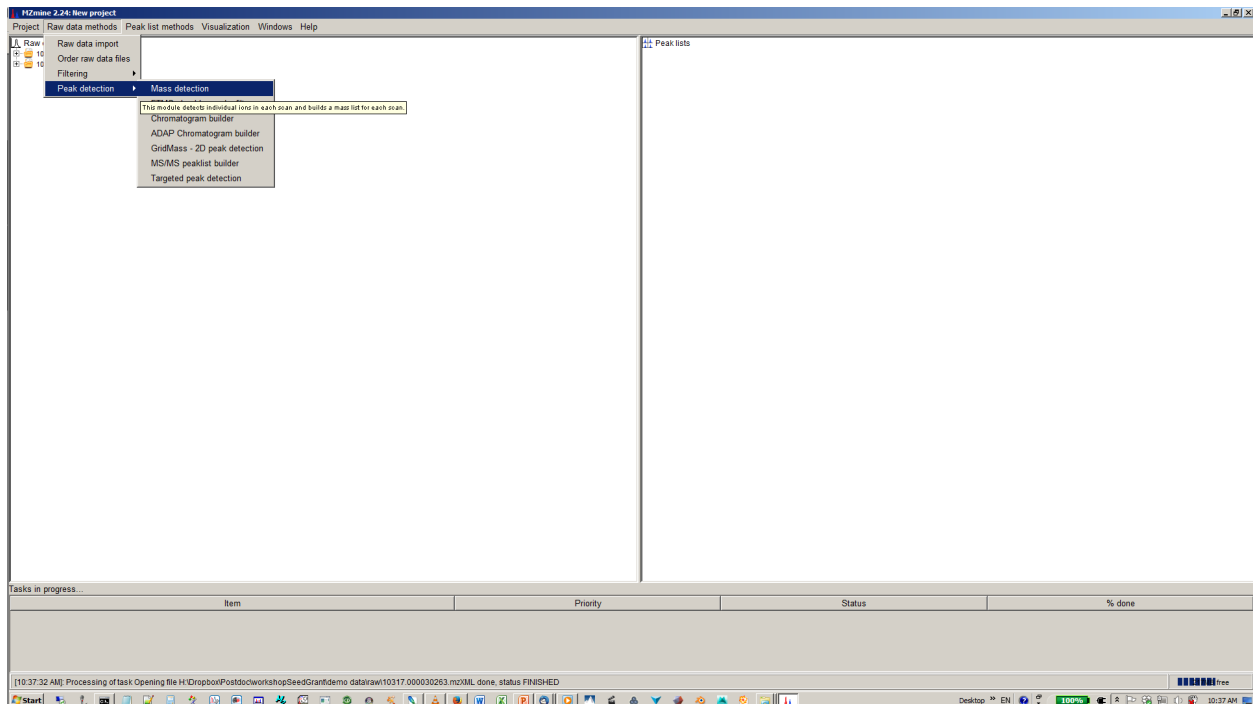
1. Start mzMine2



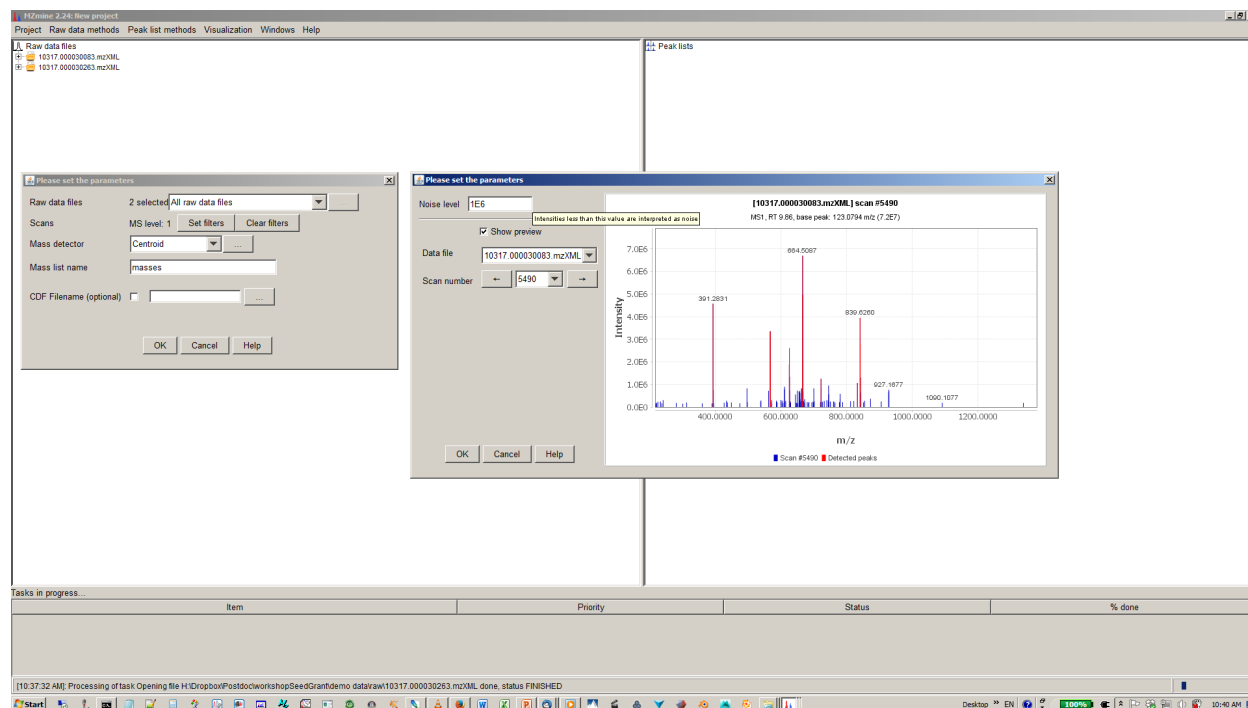
2. Click on raw data import in drop down menu and select .mzxml files



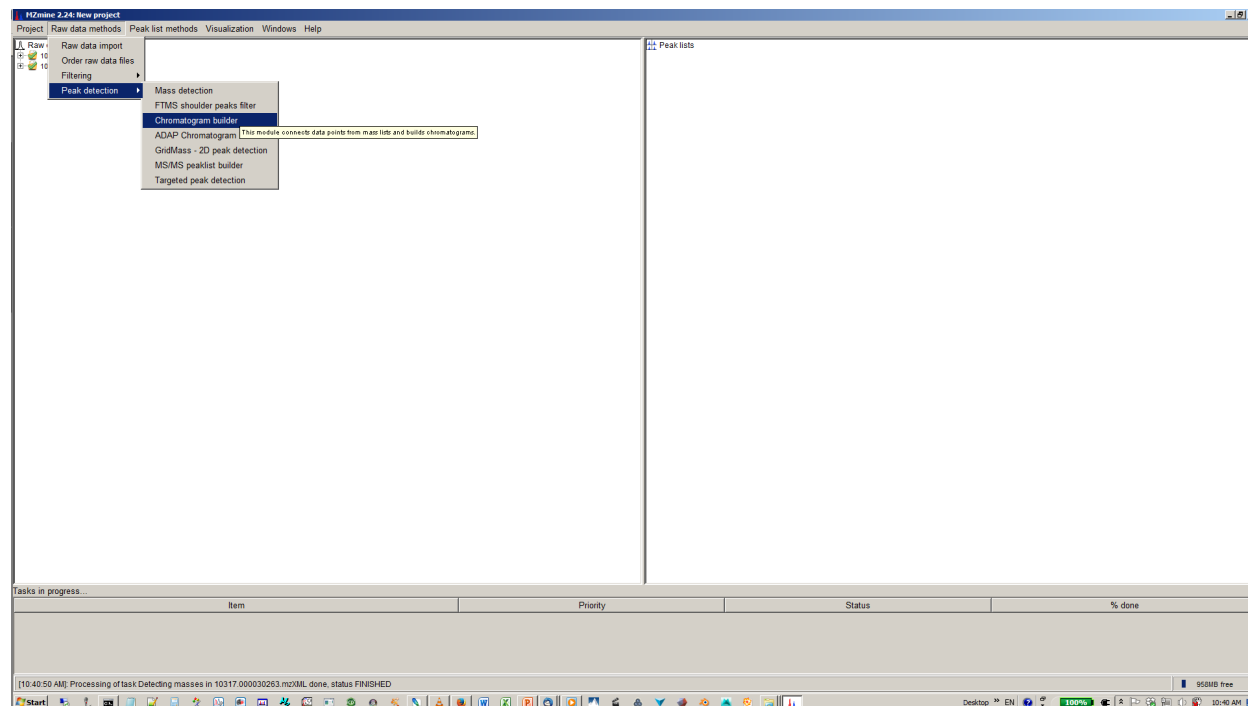
3. Click on mass detection in drop down menu



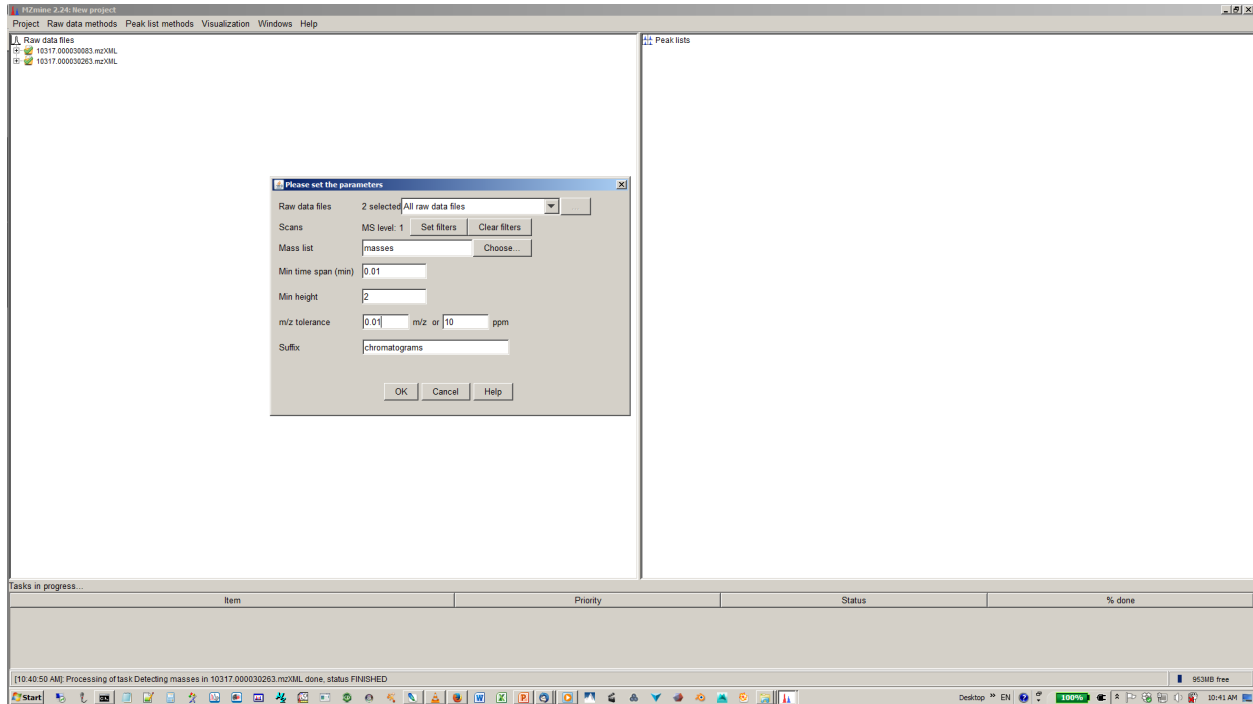
4. Specify intensity cut-off and mass list



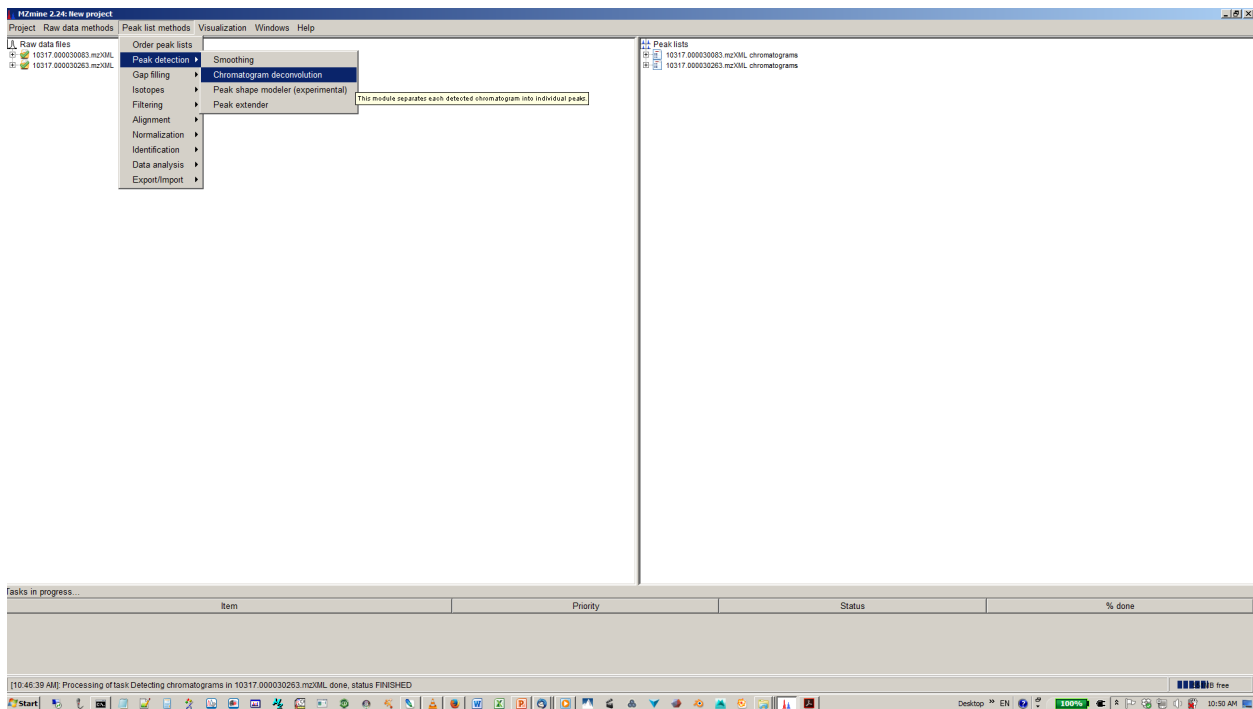
5. Build XICs with chromatogram builder



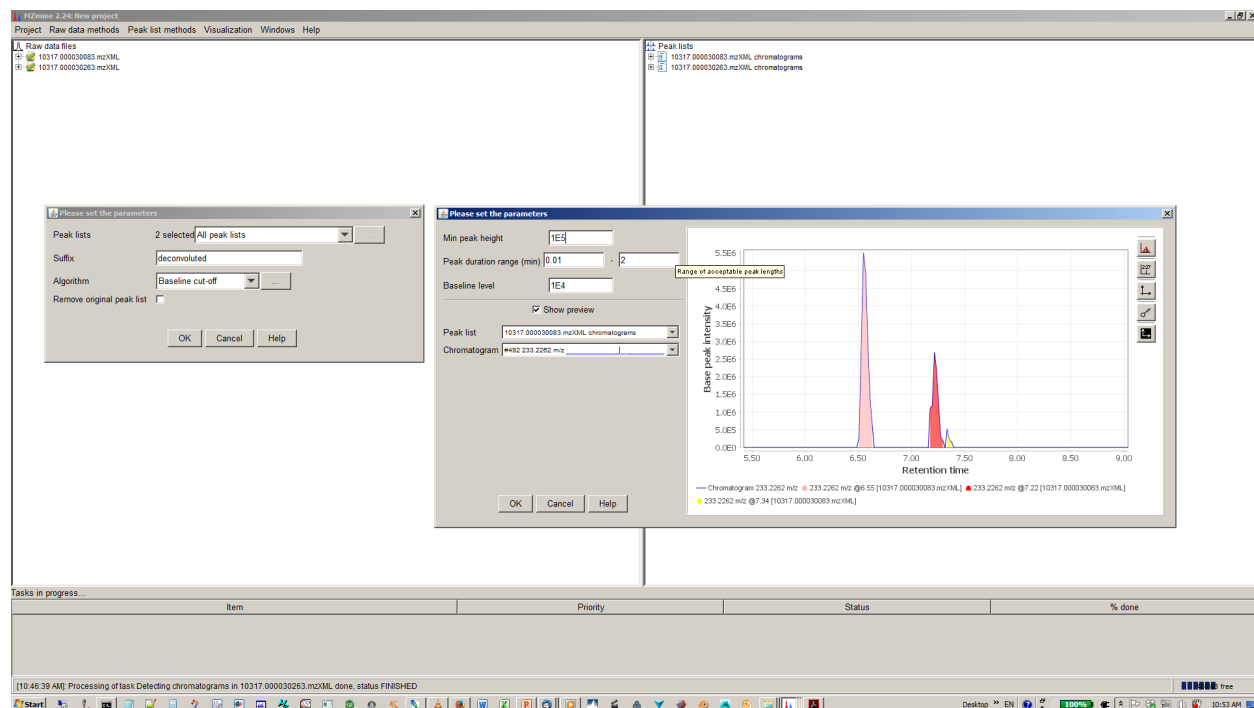
6. Specify mass list, mass tolerance min. time span and min. hight



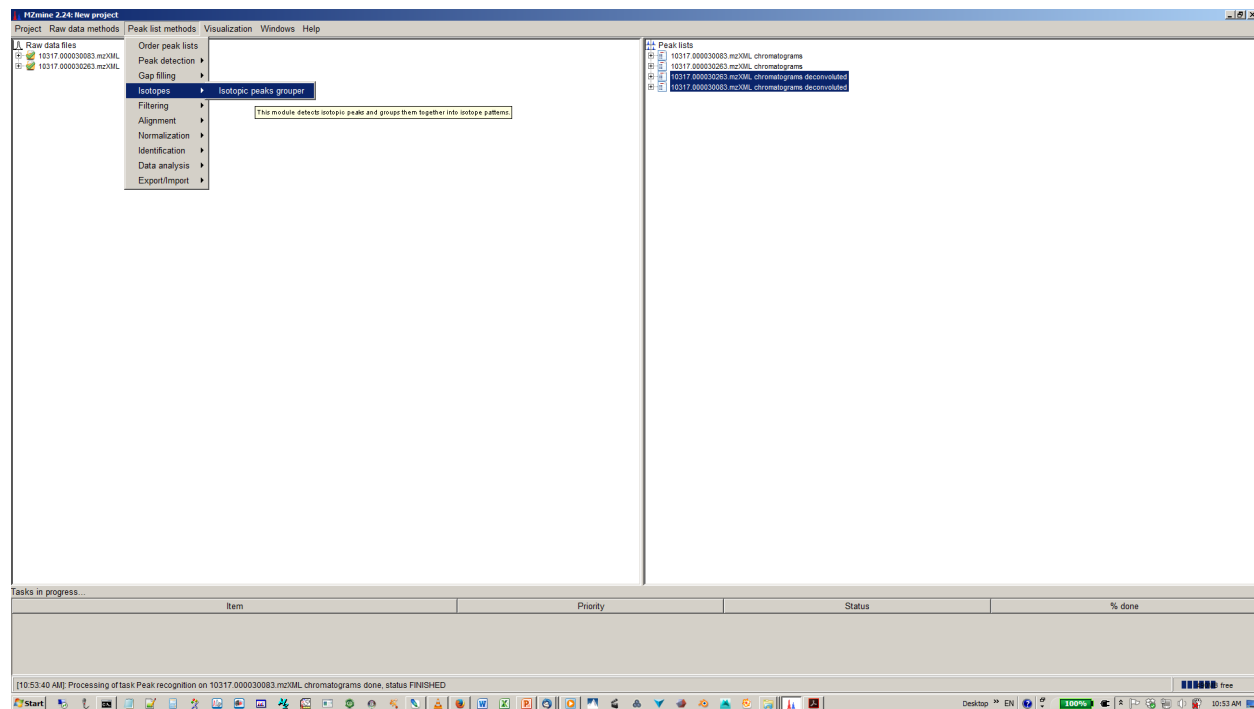
7. Deconvolute isobaric peaks with chromatogram deconvolution



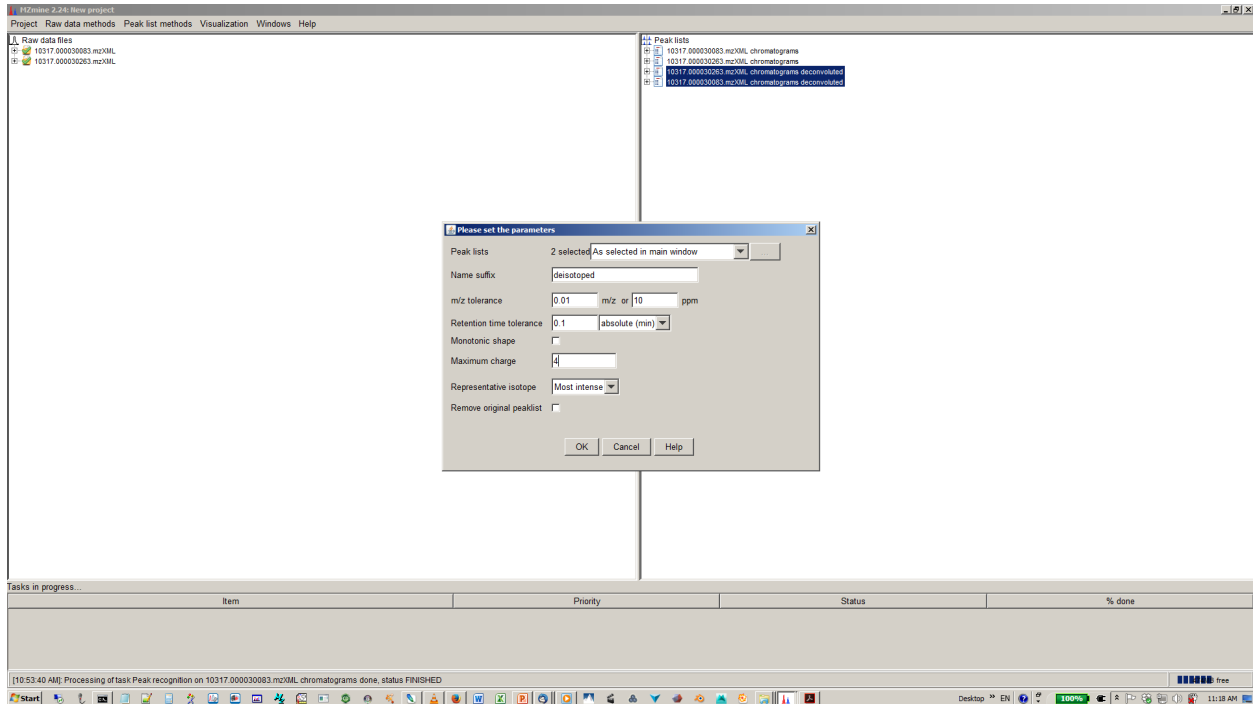
8. Specify algorithm (base line cut-off or local minimum search and parmeters)



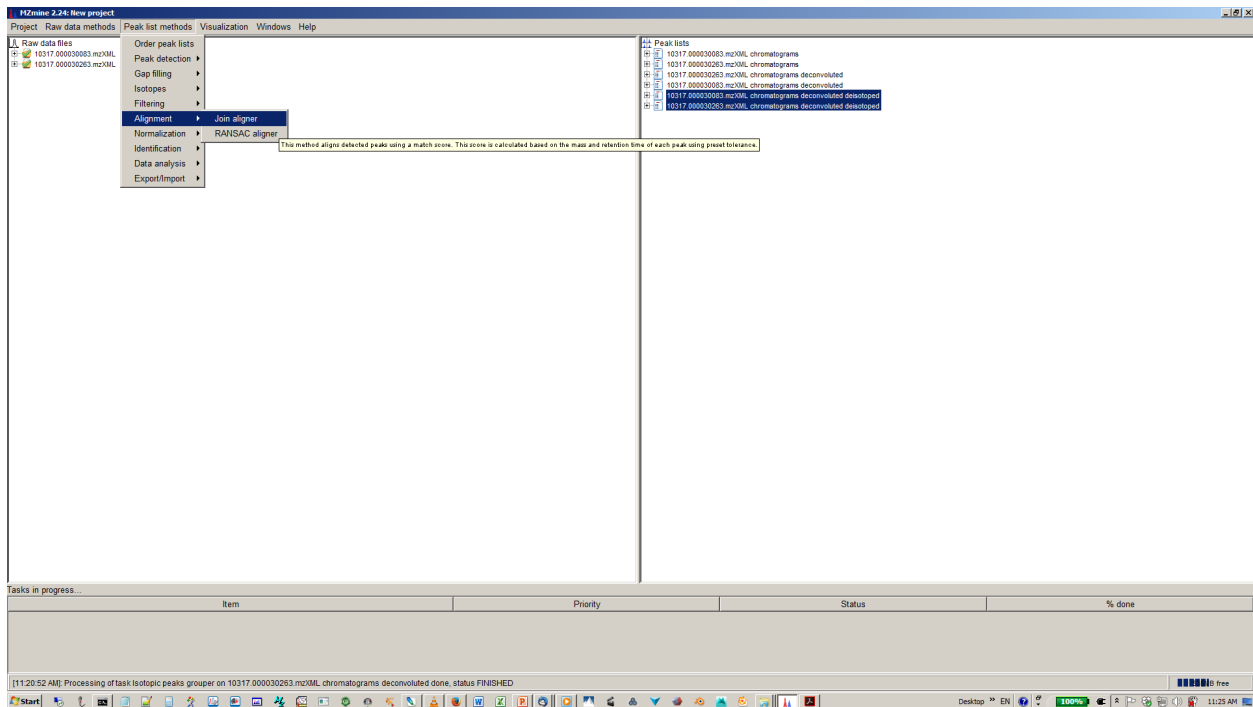
9. Perform deisotopization through isotope peak grouper



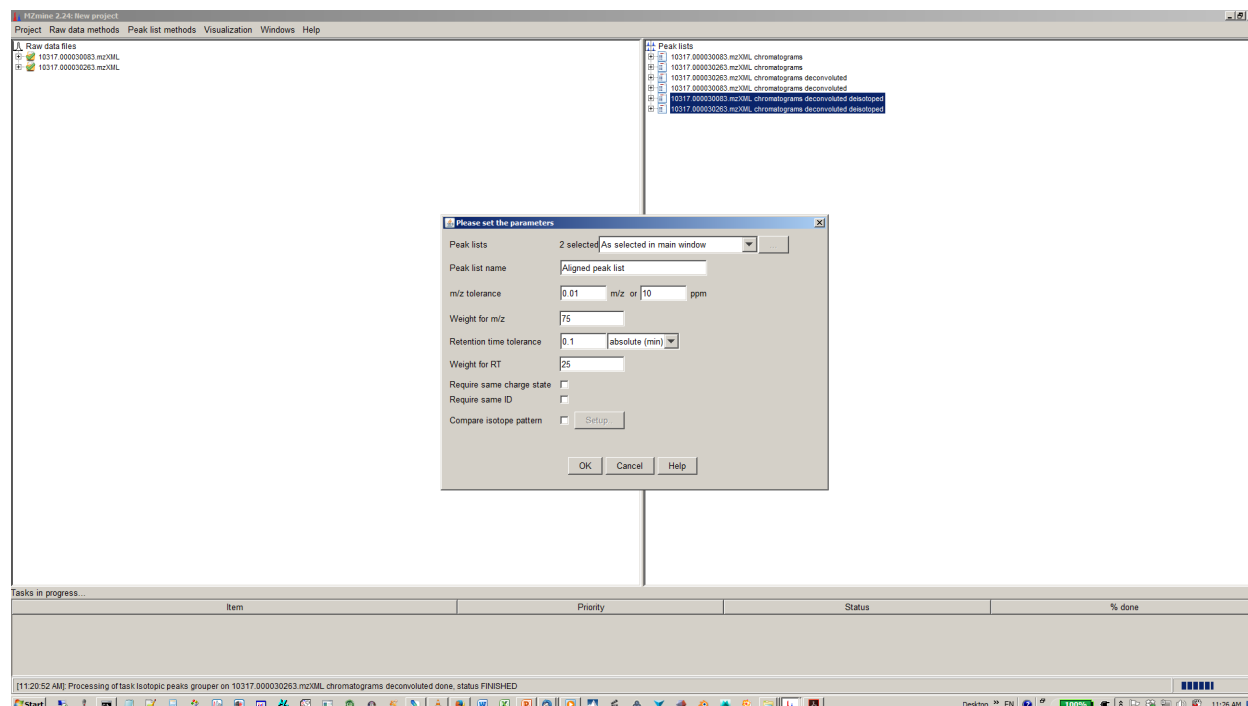
10. Specify parameters for isotope peak grouping



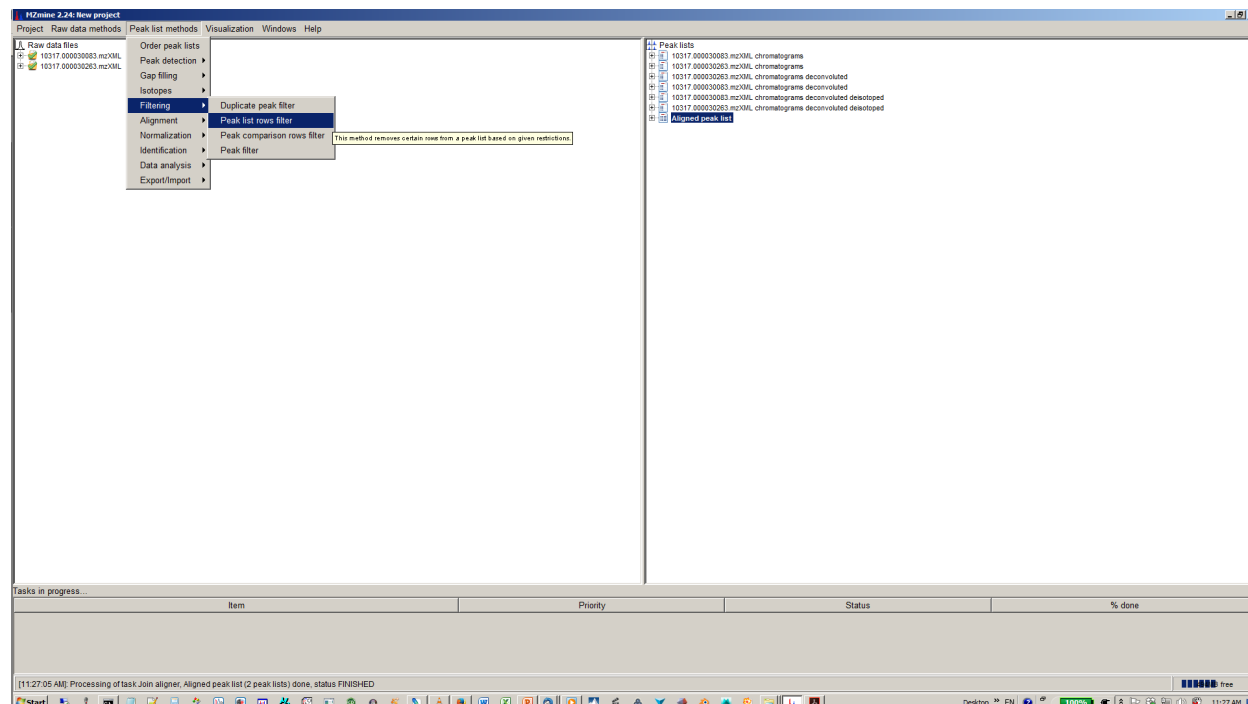
11. Align XICs from different sample to one matrix



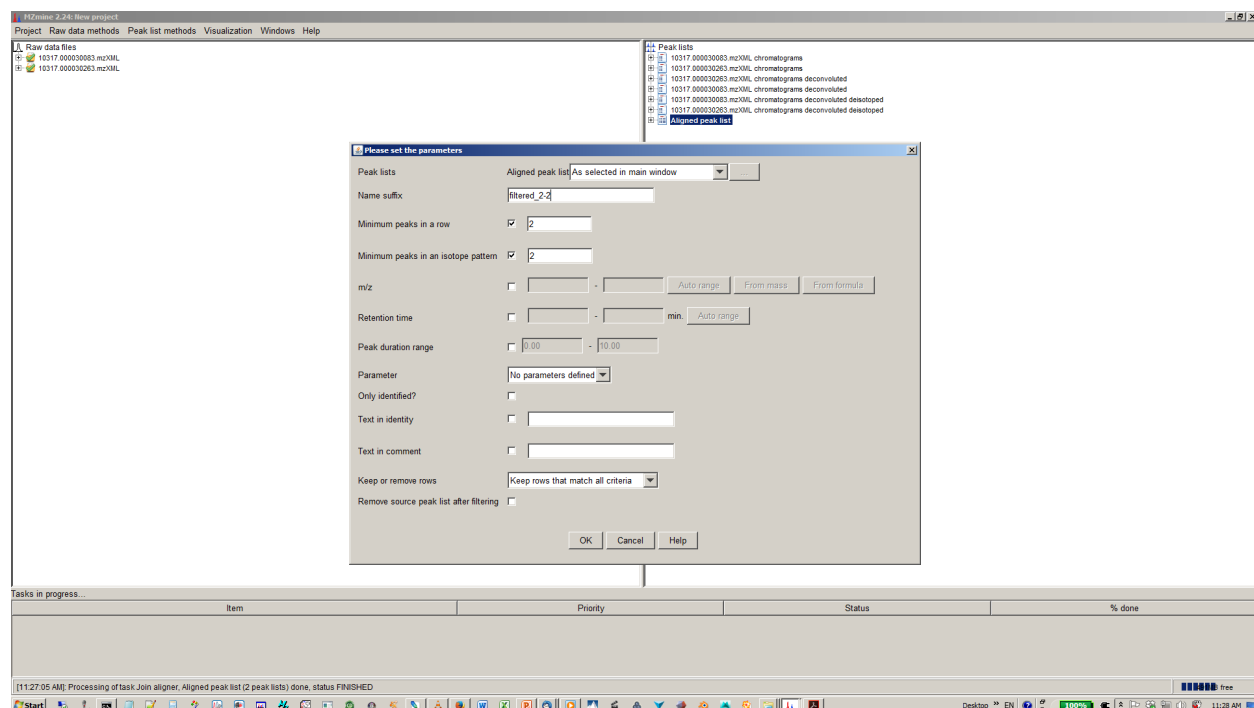
12. Specify join aligner parameters



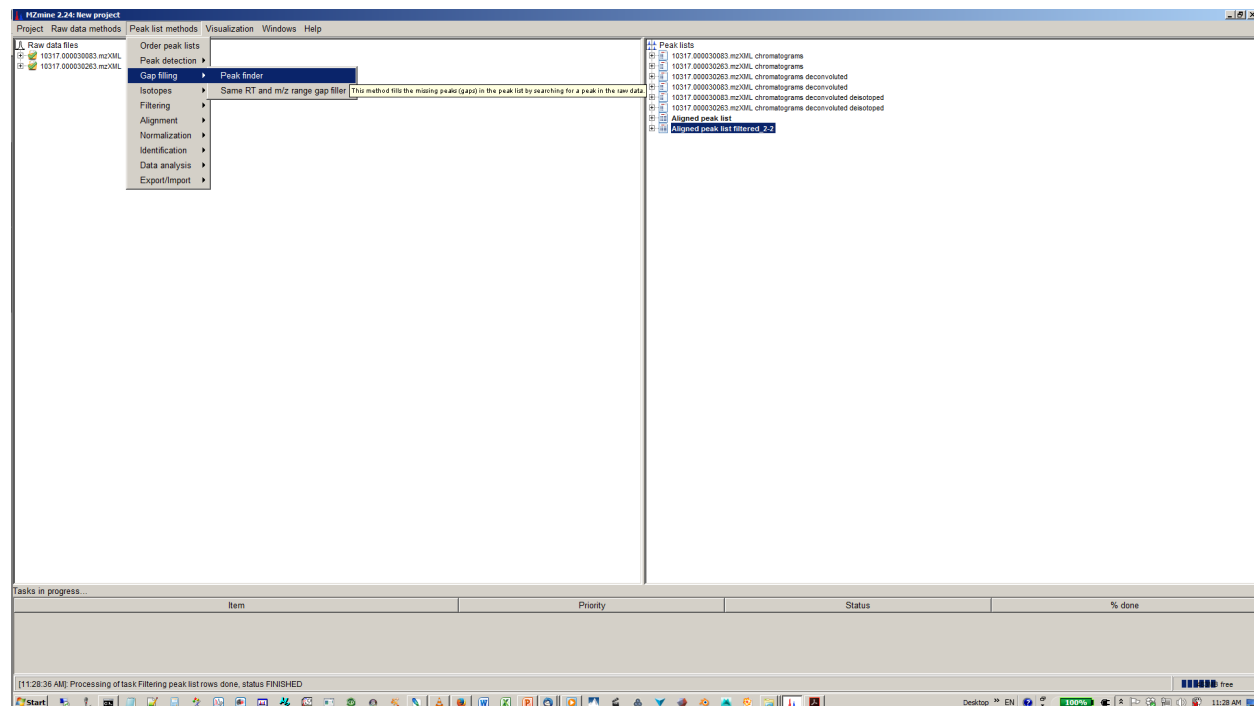
13. [optional] Filter aligned feature matrix with peak list row filter



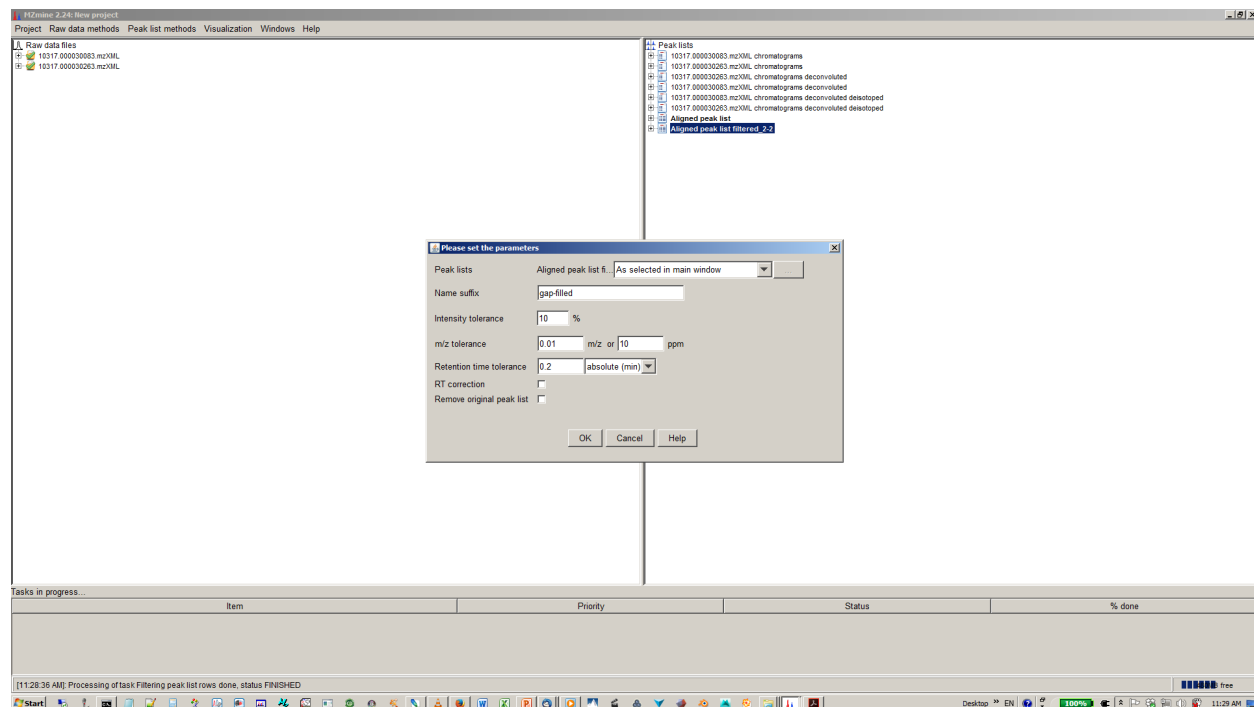
14. [optional] Depending of your experimental design use n minimum peaks in a row (n should be around the number of replicates or samples you expect to be similar) and 2-3 minimum peaks per isotope pattern



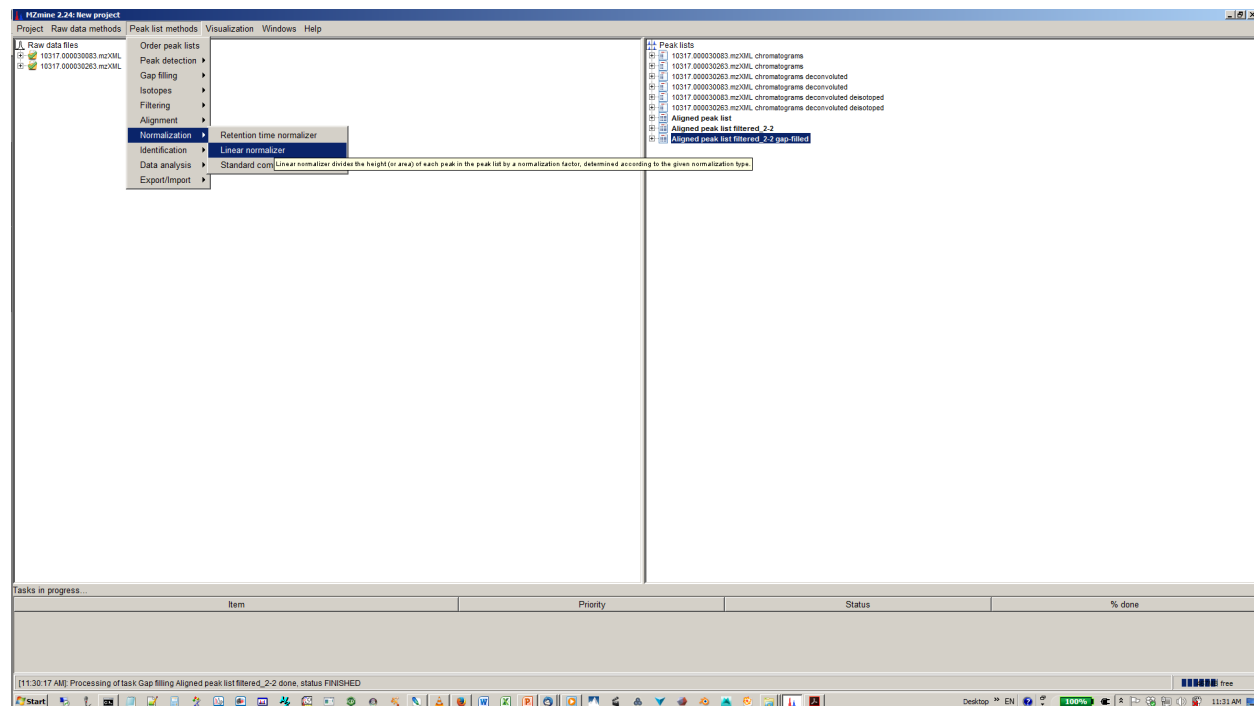
15. [optional] You gap filling the re-analyses missed peaks and fill gaps in the feature matrix



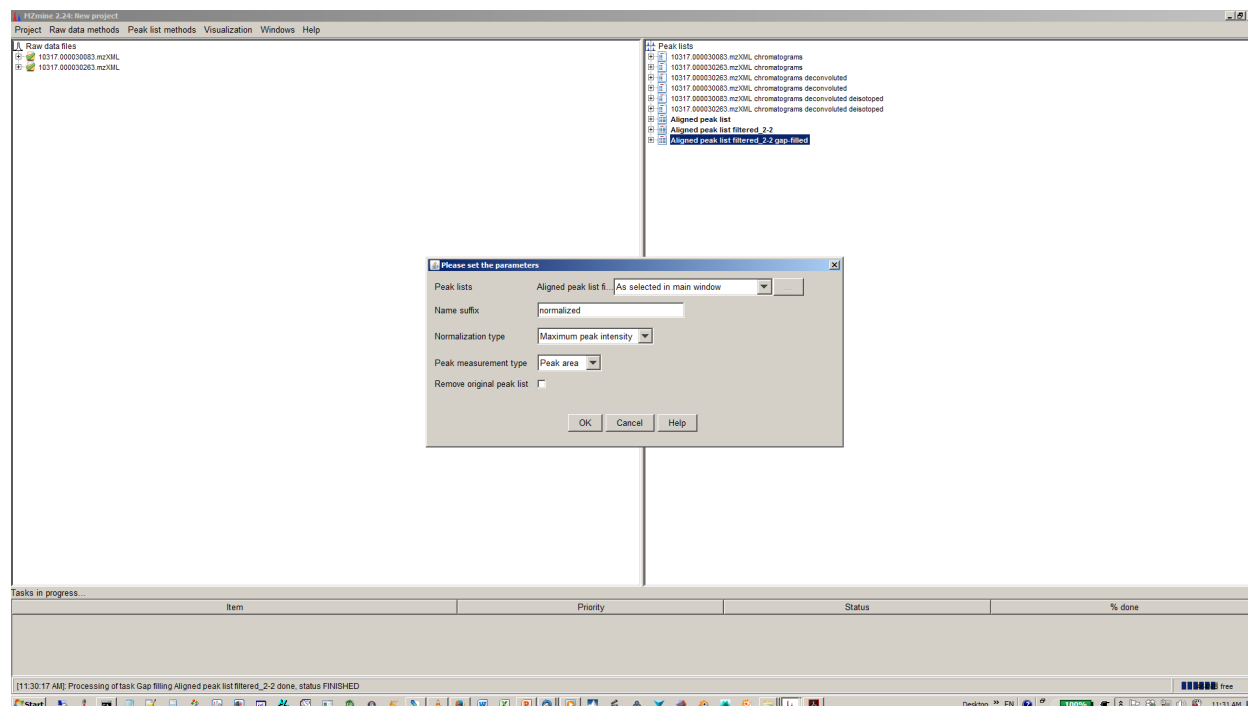
16. [optional] Depending on experimental design you can normalize your peak intensities to internal standards, TICs or total peak area.



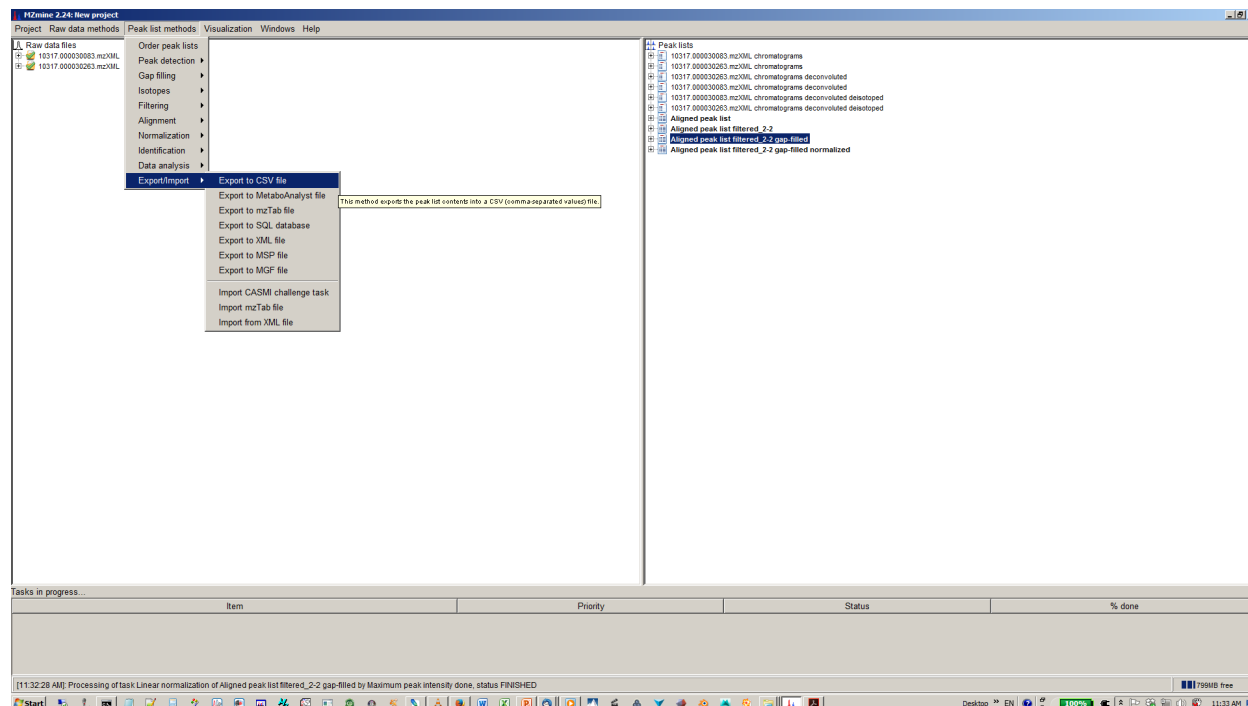
17. [optional] Specify normalization parameters

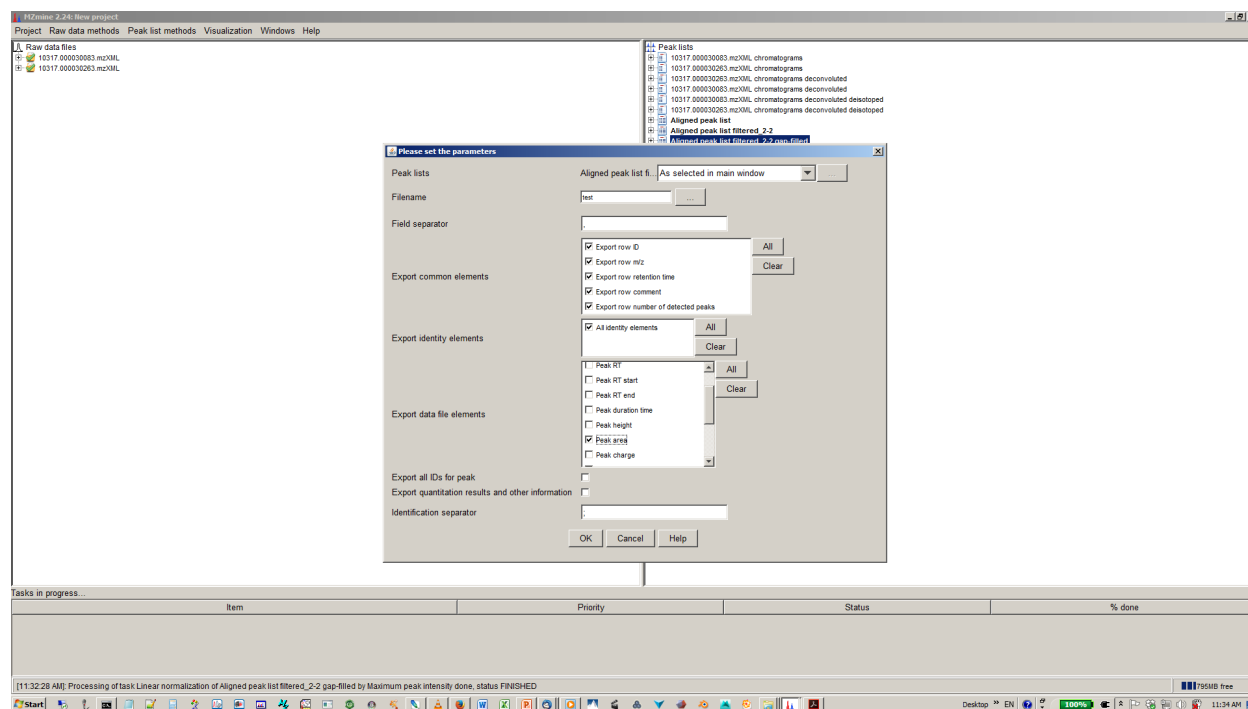


18. Export your matrix as .csv file for down stream data analysis



19. select file name and parameters you want to export





Here is also a video for [MZmine 2 documentation](#):

Metabolomics demo data in Qiita

- Refer to the Qiita documentation about Principal Coordinates Analysis (PCoA) [here](#)

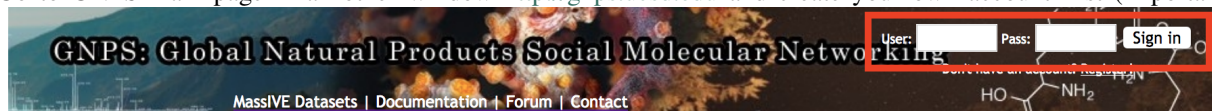
GNPS tutorial for MS/MS data annotation

Global Natural Products Social Molecular Networking [GNPS](#) web-platform provides public data set deposition and/or retrieval through the Mass Spectrometry Interactive Virtual Environment (MassIVE) data repository. The GNPS analysis infrastructure further enables online dereplication, automated molecular networking analysis, and crowdsourced MS/MS spectrum curation. Each data set added to the GNPS repository is automatically reanalyzed in the next monthly cycle of continuous identification. For more information, please check out the GNPS paper published in Nature Biotechnology by Ming et al 2016 [here](#) as well as the video and the resource on [Youtube](#), and well as on the online [documentation](#)

Tutorial: Generation of Molecular Networks in 15 minutes: Exploring MS/MS data with the GNPS Data Analysis workflow

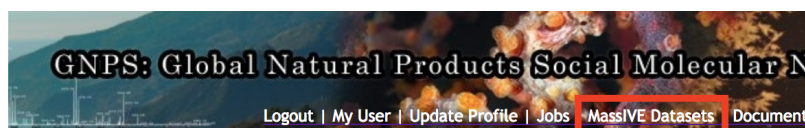
Step 1- Go to GNPS and create an account

Go to GNPS main page in an other window <http://gnps.ucsd.edu> and create your own account first (important!)



The Future of Natural Products Research and Mass Spectrometry

Step 2- Find a MS/MS dataset on MassIVE (Mass spectrometry Interactive Virtual Environment)



The Future of Natural Products Research and

A) Go to [GNPS](#) and access the MassIVE datasets repository.

B) Search for the MassIVE datasets named “GNPS Workshop” (or “GNPS_AMG_SeedGrant” for a larger example with American Gut Projects samples). Explore its content, and copy the MassIVE ID number (MSV)

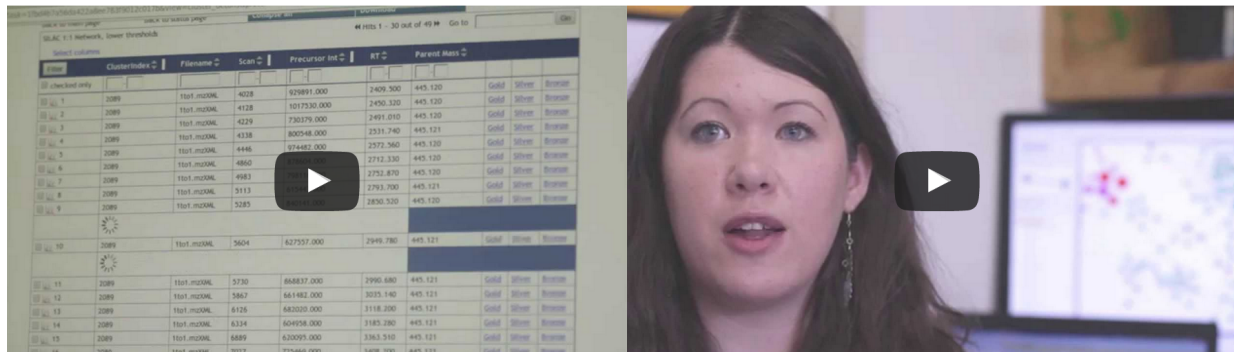
Submitted MassIVE Datasets										
Hits 1 - 4 out of 4										
Select columns										
Filter	Title	MassIVE ID	ProteomeXchange ID	Submission Type	Uploaded By	Principal Investigator	Upload Date	# Files	Total Size (KB)	
	AMG									
1	GNPS_AMG_SeedGrant	MSV000080469		Partial	rsilva	Rob Knight	Jan. 13, 2017, 10:30 AM	24	226,716	

Note: If you want to upload your own data, follow the [DorresteinLab youtube channel](#), here is the video:

Step 3 - Access to the Data Analysis workflow

Go to back [GNPS](#) main page and open the Data Analysis workflow.

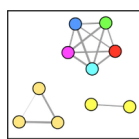
The Future of Natural Products Research and Mass Spectrometry



ClusterIndex	Filename	Scan	Precursor m/z	RT	Parent Mass
1	1001.ms.DMAL	4028	89991.000	2409.300	445.120
2	1001.ms.DMAL	4128	101720.000	2450.320	445.120
3	1001.ms.DMAL	4229	730379.000	2491.010	445.120
4	1001.ms.DMAL	4338	800948.000	2531.740	445.121
5	1001.ms.DMAL	4446	874482.000	2572.260	445.120
6	1001.ms.DMAL	4860		2712.330	445.120
7	1001.ms.DMAL	4983		2752.870	445.120
8	1001.ms.DMAL	5113		2793.700	445.121
9	1001.ms.DMAL	5285		2830.530	445.120
10	1001.ms.DMAL	5604	627557.000	2949.780	445.121
11	1001.ms.DMAL	5730	668837.000	2990.680	445.121
12	1001.ms.DMAL	5867	661482.000	3025.140	445.121
13	1001.ms.DMAL	6126	682020.000	3118.280	445.121
14	1001.ms.DMAL	6324	604958.000	3185.280	445.121
15	1001.ms.DMAL	6889	630093.000	3363.510	445.121
16	1001.ms.DMAL	7072	723469.000	3408.380	445.121

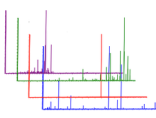
[Tweet](#) [Share](#)

Data Analysis



The [Data Analysis](#) portal will allow you to organize and visualize your mass spectrometry data. Leveraging the molecular networking techniques, there are additional tools to aid in understanding the unknowns in your sample. Check out the [documentation](#) and live [demo](#). Further, a separate [dereplication workflow](#) is provided as a standalone workflow.

Create Public Massive Datasets



[Submit](#) your own data to be made public Massive datasets. These Massive datasets must be **prefixed with GNPS** to be visible to other GNPS users. Take advantage of [continuous identification](#) to learn more about your dataset after publication automatically. New hits to the community curated libraries and related datasets are reported. [Documentation](#)

Step 4 - Configure and launch the Data Analysis workflow

Workflow Selection

Title:

Search Protocol:

Networking Parameter Presets

Basic Options

Spectral Library: 0 files and 1 folder are selected [To import libraries for search click here](#)

Spectrum Files (Required): 0 files and 1 folder are selected [See here for further documentation about molecular networking.](#)

Spectrum Files G2: [Click Here here to run a demo molecular network.](#)

Spectrum Files G3:

Spectrum Files G4:

Spectrum Files G5:

Spectrum Files G6: [For custom group/attribute documentation click here](#)

Precursor Ion Mass Tolerance: Da Fragment Ion Mass Tolerance: Da

Advanced Network Options

Advanced Library Search Options

Advanced Filtering Options

Advanced Output Options

Workflow Submission

Email me at

A) Indicate a Title.

Select Input Files

Select Input Files

Library Files

Spectrum Files G1

Spectrum Files G2

Spectrum Files G3

Spectrum Files G4

Spectrum Files G5

Spectrum Files G6

Group Mapping

Attribute Mapping

Selected Files

Selected Library Files

Selected Spectrum Files

Selected Spectrum Files

Selected Spectrum Files

Selected Spectrum Files

Selected Spectrum Files

Selected Spectrum Files

Selected Spectrum Files

Selected Group Mapping

B) Click on Spectrum Files (required)

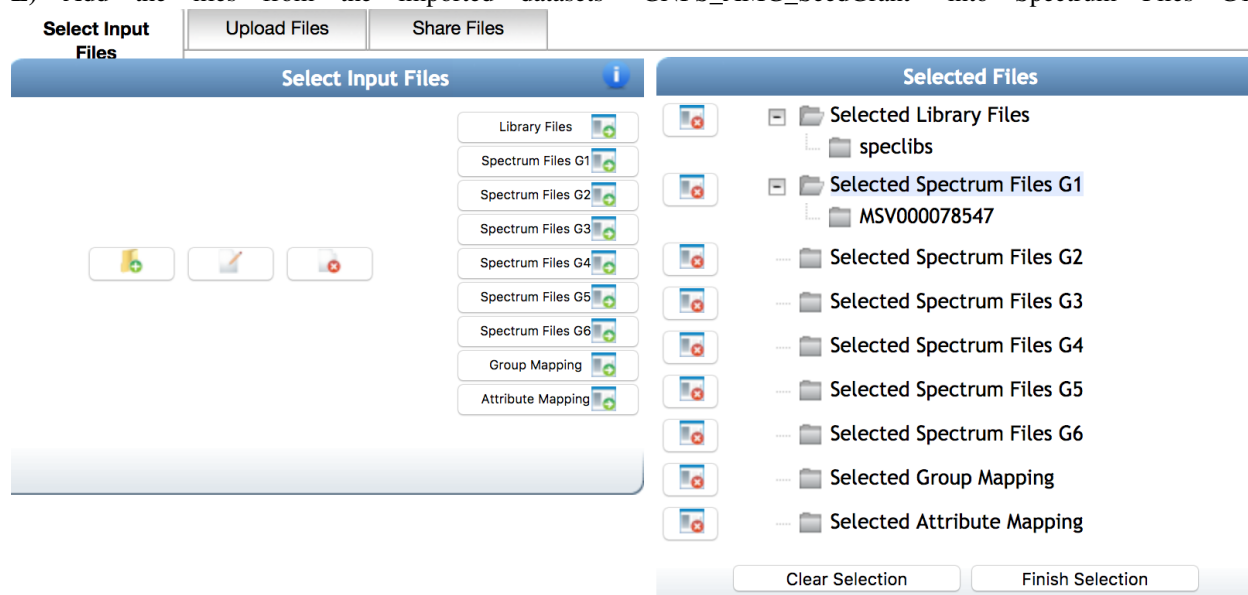
C) Go to the Share Files spreadsheet and import the Massive dataset files for the “GNPS workshop” or

figs/GNPS_import.png

“GNPS_AMG_SeedGrant” with the Import Data Share (use the MassIVE ID).

D) Go back to the Select Input Files spreadsheet.

E) Add the files from the imported datasets “GNPS_AMG_SeedGrant” into Spectrum Files G1.



F) Validate the selection with Finish Selection button.

G) Modify parameters to meet high-resolution mass spectrometry: Precursor Ion Mass Tolerance (0.02), Fragment Ion Mass Tolerance (0.02), Min Pairs Cos (0.6), Minimum Matched Fragment Ions (2), Minimum cluster size (use 1)

Workflow Selection

Title: Workshop AMG

Search Protocol: None Reset Form Save as Protocol

Networking Parameter Presets

Small Data Preset

Medium Data Preset

Big Data Preset

Basic Options

Spectral Library: Select Input Files 0 files and 1 folder are selected To import libraries for search click [here](#)

Spectrum Files (Required): Select Input Files 0 files and 1 folder are selected See [here](#) for further documentation about molecular networking.

Spectrum Files G2: Select Input Files Click Here [here](#) to run a demo molecular network.

Spectrum Files G3: Select Input Files

Spectrum Files G4: Select Input Files

Spectrum Files G5: Select Input Files

Spectrum Files G6: Select Input Files For custom group/attribute documentation click [here](#)

Precursor Ion Mass Tolerance: 2.0 Da Fragment Ion Mass Tolerance: 0.5 Da

Advanced Network Options

Show Fields

Advanced Library Search Options

Show Fields

Advanced Filtering Options

Show Fields

Advanced Output Options

Show Fields

Workflow Submission

Email me at lnothiasscaglia@ucsd.edu

Submit

H) Launch the Data Analysis workflow using the Submit button.

Step 5 - Visualize the Data Analysis workflow output

A) Return to GNPS main page and go to the Jobs page. Please find here an example of GNPS data analysis output with American Gut Project.

The Future of Natural Products Research and Mass Spectrometry

1.11. GNPS tutorial for MS/MS data annotation

51

Job Status

Workflow

METABOLOMICS-SNETS

Status

DONE

[Clone]

[View All Library Hits](#)

[View All Clusters With IDs](#)

[View All Compounds](#)

[Restart][Delete]

Methods and Citation for Manuscripts

[Networking Parameters and Written Network Description]

Experimental Views

[Reanalyze Cluster Spectra | View Raw Spectra | Topology Signatures | Topology Signatures Histogram]

Auxiliary Views

[View Network, Node Centric | View Network Pairs | Networking Statistics]

Advanced Views - Networking Graphs

[Nodes, MZ Histogram | Edges, MZ Delta Histogram | Edges, Score vs MZ Delta Plot | Library Search, PPM Error Histogram]

Community Matches

[Dataset Matches]

Network Visualizations

[View Spectral Families (In Browser Network Visualizer) | Network Summarizing Graphs]

Export

[Download Clustered Data | Download Cytoscape Data | Download Bucket Table | Make Public Dataset]

User

lfnothias (lfnothiascaglia@ucsd.edu), UCSD, Dorrestein Lab

Title

Workshop AMG

Re-Analyze Task Outputs

[Import to Re-analyze Task Data](#) [Attach Reanalysis Results to Dataset](#)

Date Created

2017-01-19 15:39:19.0

Execution Time

9 minutes 12 seconds

B) Explore the molecule annotated using public spectral library available on GNPS. Click on [View All Library Hits](#).

Workshop AMG

Hits 1 - 30 out of 43

Go to

Go

Filter

ViewLib

Compound Name

Library Class

Cosine

MZErrorPPM

MassDiff

LibMZ

Instrument

IonMode

PI

Ion Source

1

[ViewLib](#)

Abyssomin B

Gold

0.95

2657

1.005

378.16

Hybrid FT

Positive

Dorrestein

DI-ESI

Query

Library

Back to main page

[Back to status page](#)

Collapse all

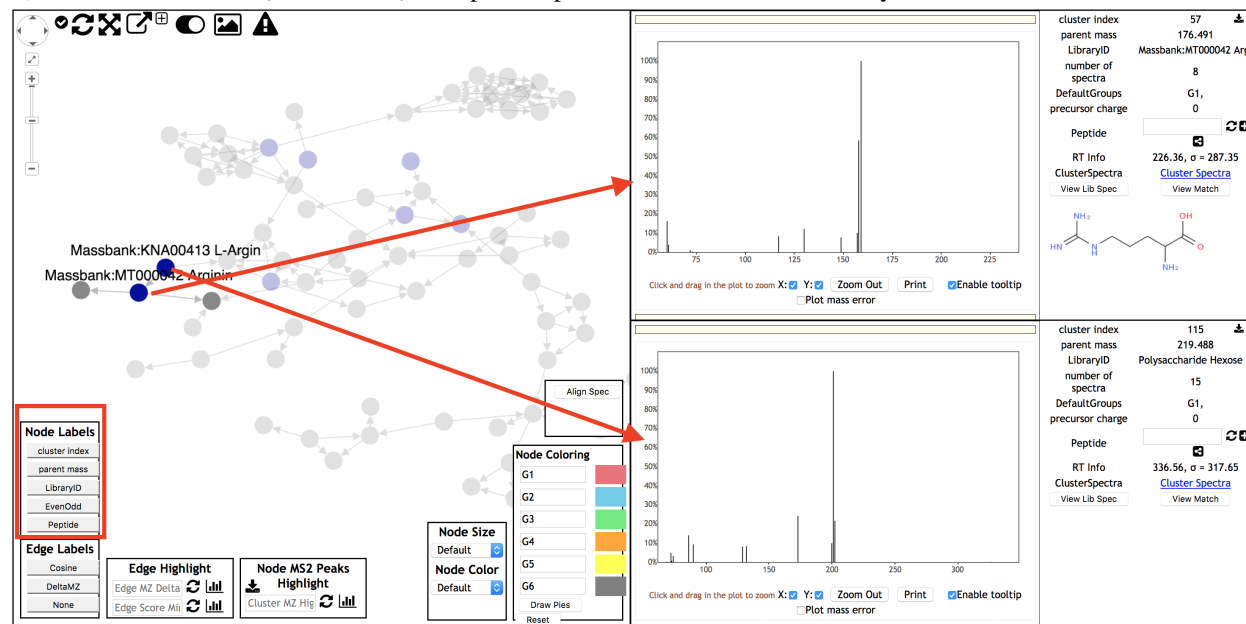
Download

C) Go back to the Status Page

D) Click on the [View Raw](#) Spectral families and visualize the molecular network

Status	Auxiliary Views [View Network , Node Centric View Network Pairs Networking Statistics]
	Advanced Views - Networking Graphs [Nodes, MZ Histogram Edges, MZ Delta Histogram Edges, Score vs MZ Delta Plot Library Search, PPM Error Histogram]
	Community Matches [Dataset Matches]
	Network Visualizations [View Spectral Families (In Browser Network Visualizer) Network Summarizing Graphs]
	Export [Download Clustered Data Download Cytoscape Data Download Bucket Table Make Public Dataset]
User	lfnothias (lfnothiascaglia@ucsd.edu), UCSD, Dorrestein Lab
Title	Workshop AMG

E) In Node Labels (bottom left), map the parent mass, or the LibraryID, in the molecular network.



F) Visualize a first MS/MS spectrum by left-clicking on one node. Visualize a second MS/MS spectrum by right-clicking on a second node.

More on navigating into the results with the following video: