

# Protein MS isn't the only MS...

## Programming in R for Metabolomics

### Mass Spectrometry

**Chase Clark**  
**PhD Candidate**  
**Murphy Lab**

**Dept. of Pharmaceutical Sciences**  
**Center for Biomolecular Sciences**  
**University of Illinois at Chicago**



@ChasingMicrobes  
chasemc.github.io

May Institute Future Developers Meeting



~~Protein MS isn't the only MS...~~

~~Programming in R for Metabolomics~~

~~Mass Spectrometry~~

Programming in R for Mass-omics

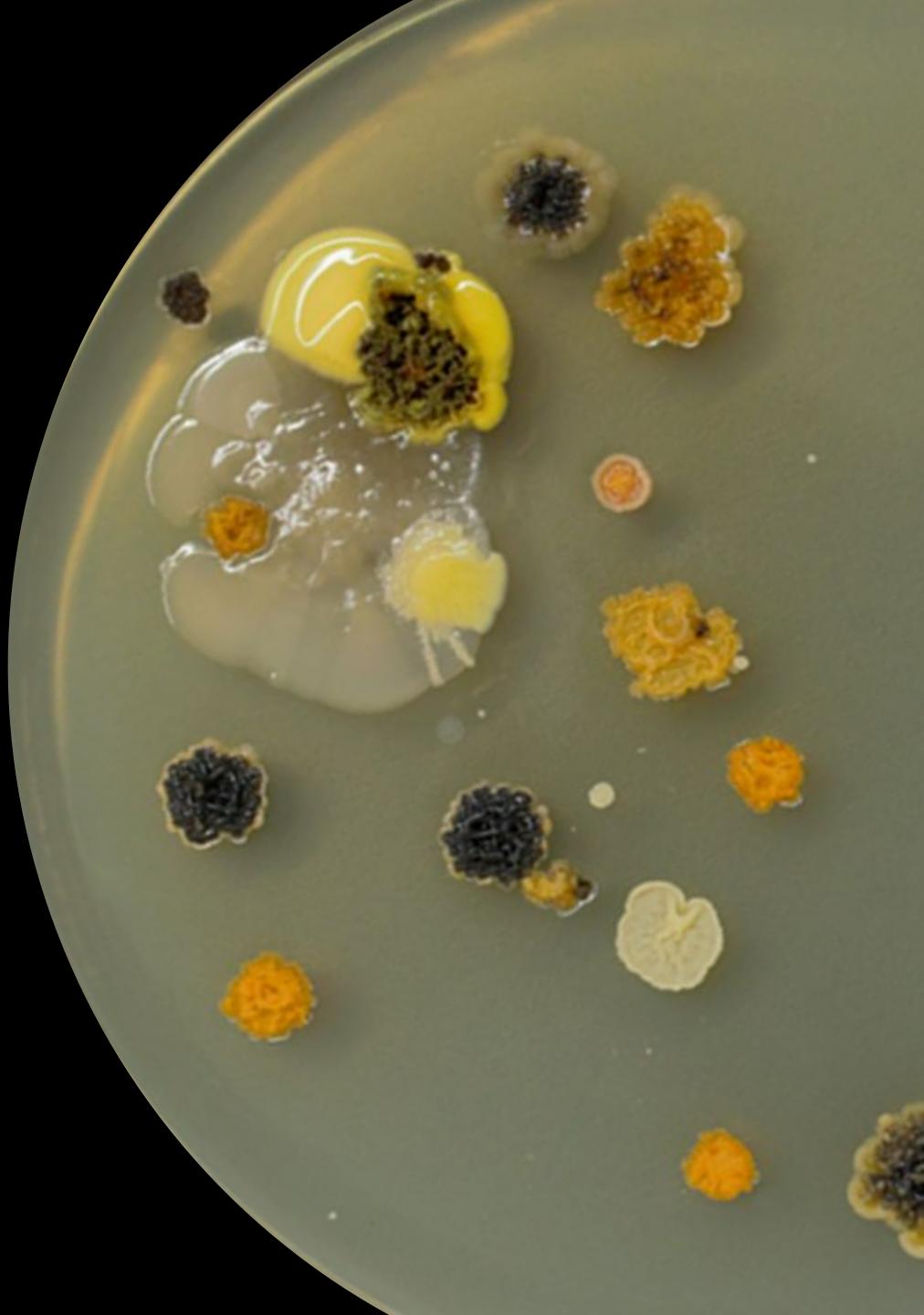
**Chase Clark**  
**PhD Candidate**  
**Murphy Lab**

**Dept. of Pharmaceutical Sciences**  
**Center for Biomolecular Sciences**  
**University of Illinois at Chicago**



@ChasingMicrobes  
chasemc.github.io

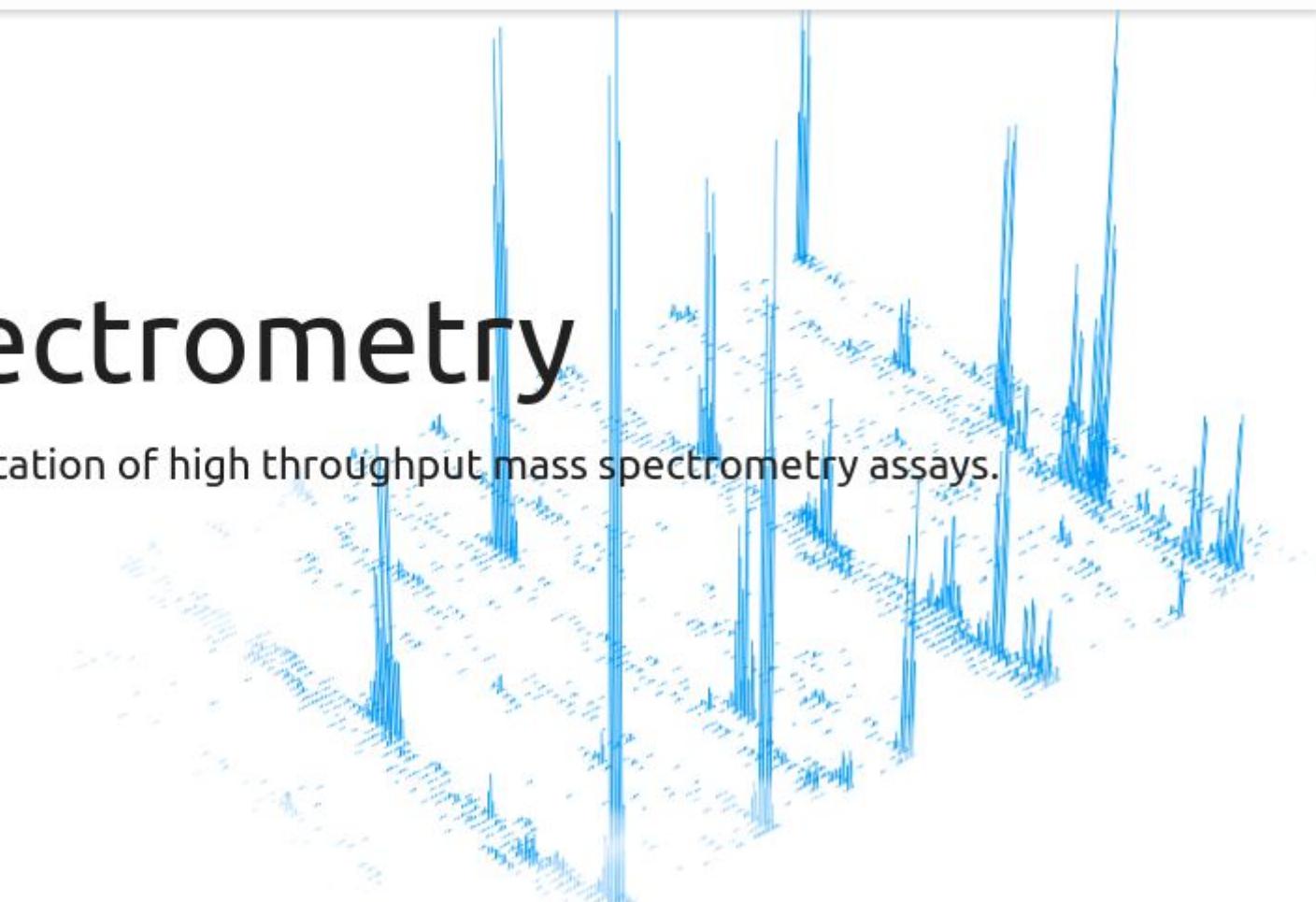
May Institute Future Developers Meeting



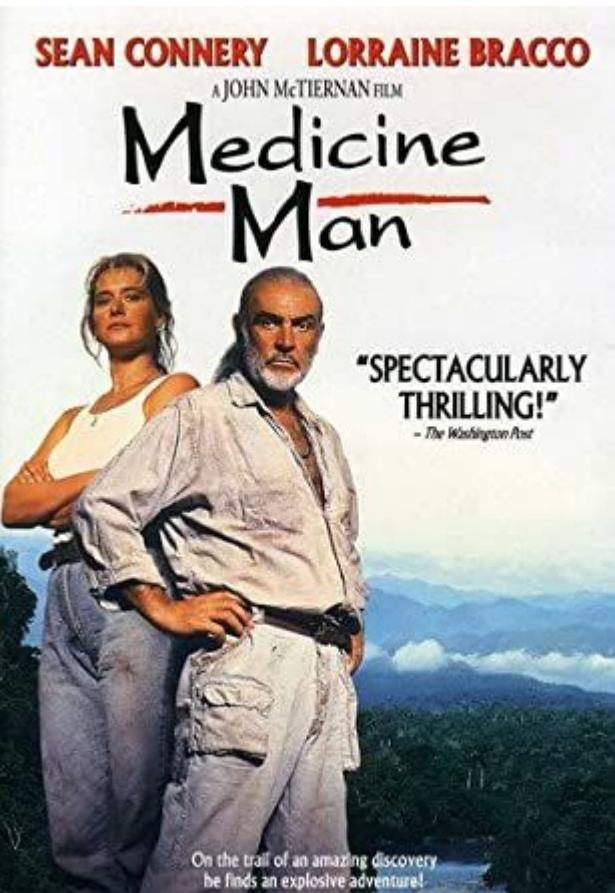
[RforMassSpectrometry](#)<https://www.rformassspectrometry.org>[Home](#)[Packages](#)[Contact](#)

# R for Mass Spectrometry

R software for the analysis and interpretation of high throughput mass spectrometry assays.

[Contact](#)

# Pharmacognosy



# Pharmacognosy

The study of the physical, chemical, biochemical and biological properties of drugs, drug substances or potential drugs of natural origin as well as the search for new drugs from natural sources.



SCIENCE MUSEUM Why scientists are hunting for bacteria in Iceland

# FROM THE LAB

SCIENCE MUSEUM

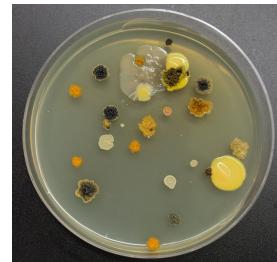
Watch later Share

A video thumbnail for a Science Museum video. The title is "Why scientists are hunting for bacteria in Iceland". The image shows two scientists in a boat on choppy water, with a snowy, volcanic landscape in the background. A play button icon is in the center. There are also "Watch later" and "Share" buttons in the top right corner.

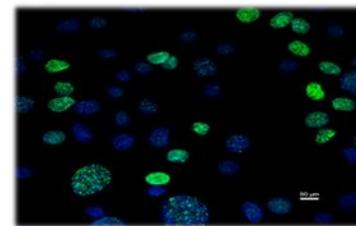
# Discovery Pipeline (Simplified)



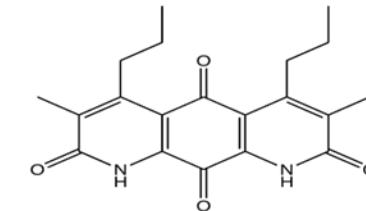
Collect  
Samples



Grow  
Bacteria



Antibiotic  
Activity?



Elucidate  
Antibiotic



Drug

# Which bacteria do we study?

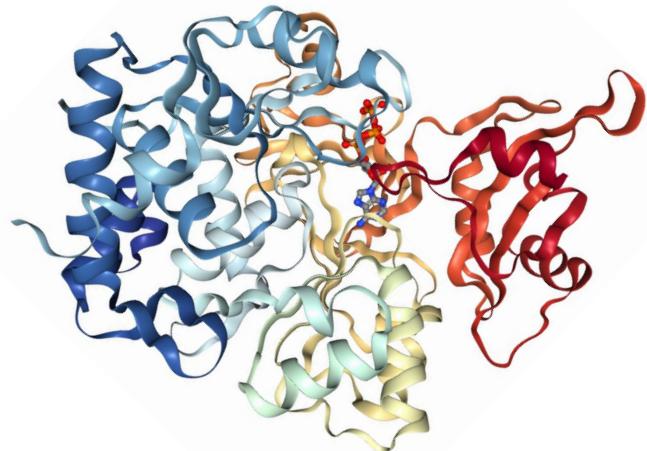


Dr. Laura Sanchez

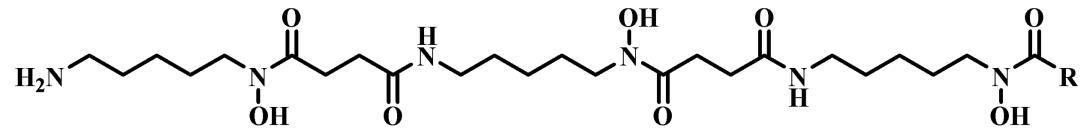


# Collect Two Fingerprints Per Colony

Protein Spectra  
2-20 kDa  
Pseudo-Phylogenetic Groupings



Small Molecule Spectra  
0.2-2 kDa  
Measuring Produced Metabolites



# Domain Knowledge Existed in R



## Species Identification using MALDIquant

Sebastian Gibb\* and Korbinian Strimmer †

June 8, 2015

### Abstract

This vignette describes how to use MALDIquant for species identification.

# But I didn't know to code!



# IDBac

## Contents

### Sample Handling

- [Cell Culture](#) pg. 1
- [MALDI Sample Preparation](#) pg. 2
- [MALDI Data Acquisition](#)

### MALDI

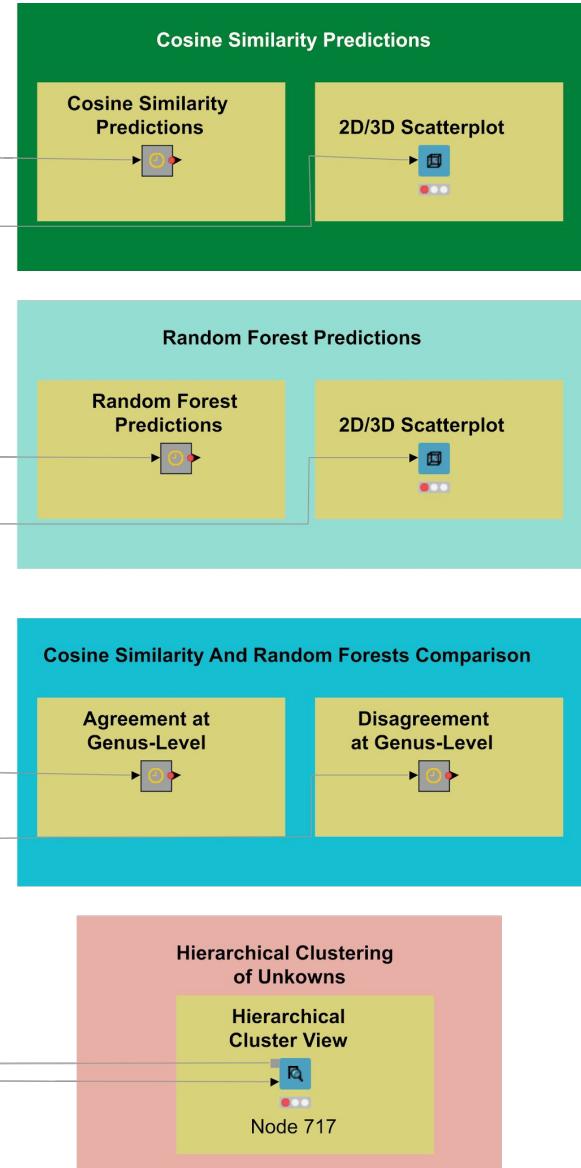
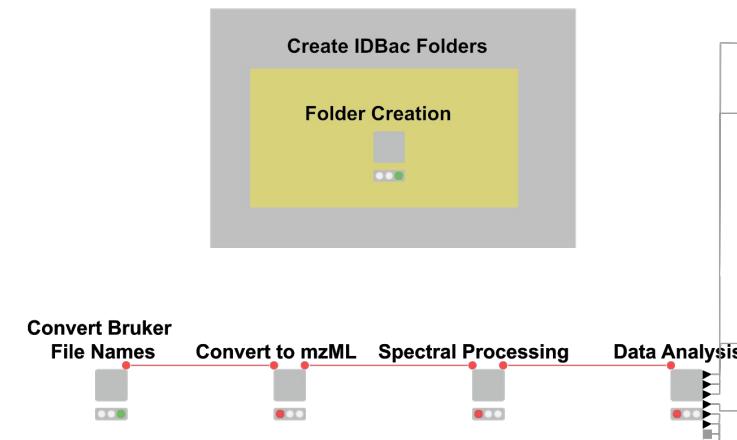
- [Instrument Configuration](#)
- [Data-File Handling](#)
- [Sample Labelling](#)
- [Software for Data Analysis](#)
  - [Installations Overview](#)
  - [Installations: Step by Step](#)
  - [R, KNIME, GEPHI, msconvert](#)
  - [IDBac Setup](#)

### Analyzing Your Unknowns

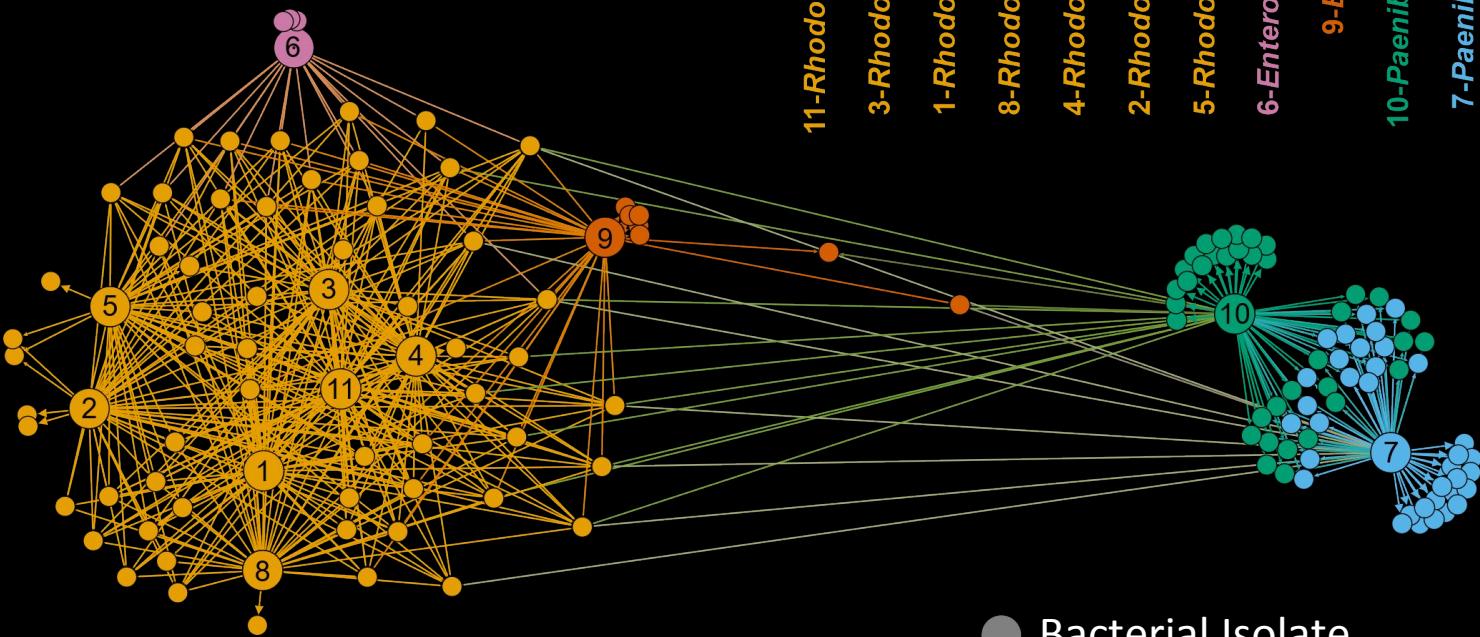
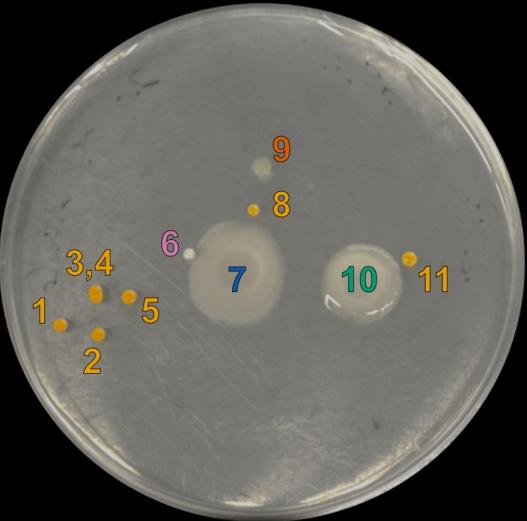
- [Characterization by protein analysis](#)
- [Small-molecule networks](#)

### Appendix A: A1 Media Recipe

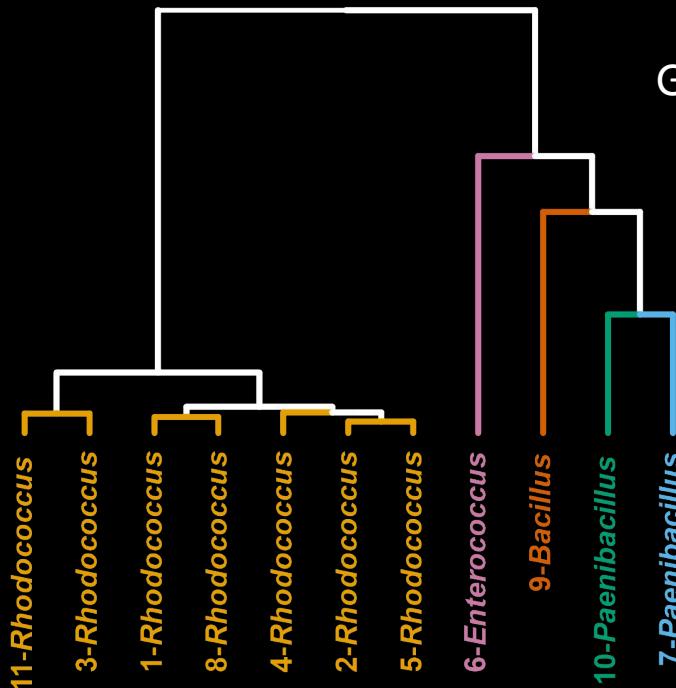
### Appendix B: Automated Bruker File Rename



# IDBac



Visualized small molecule data



Grouped using protein spectra

- Bacterial Isolate
- $m/z$  Peak

```

Package: IDBacApp
Type: Package
Title: MALDI-TOF MS Protein and Small Molecule Analysis
Version: 1.2.0.9000
Author: Chase Clark
Maintainer: Chase Clark
Description: This app allows users to take raw data MALDI-TOF MS files
and create bundled mzXML files for easy sharing/database creation as well as data analysis
(hierarchical clustering, PCA, Spectra Comparison, Metabolite Association Network).
License: GPL-3
Encoding: UTF-8
LazyData: true
biocViews:
Imports:
  ape (>= 5.3),
  colourpicker (>= 1.0),
  coop (>= 0.6-2),
  data.table (>= 1.12.2),
  DBI (>= 1.0.0),
  digest (>= 0.6.20),
  fst (>= 0.9.0),
  ggplot2 (>= 3.2.0),
  glue (>= 1.3.1),
  graphics,
  grDevices,
  httr (>= 1.4.1),
  igraph (>= 1.2.4.1),
  irrlba (>= 2.3.3),
  jsonlite (>= 1.6),
  magrittr (>= 1.5),
  MALDIquant (>= 1.19.3),
  MALDIquantForeign (>= 0.12),
  mzR,
  networkD3 (>= 0.4),
  plotly (>= 4.9.0),
  pool (>= 0.1.4.2),
  remotes,
  reshape2 (>= 1.4.3),
  rhandsonable (>= 0.3.7),
  rmarkdown (>= 1.14),
  RSQLite (>= 2.1.2),
  Rtsne (>= 0.15),
  S4Vectors,
  shinyCSSloaders (>= 0.2.0),
  sigmajs (>= 0.1.3),
  stats,
  svglite (>= 1.2.2),
  utils,
  Matrix,
  dendextend (>= 1.12.0),
  shiny (>= 1.3.2),
  curl
Suggests:
  knitr (>= 1.23),
  testthat (>= 2.2.1)
RoxygenNote: 7.1.0
VignetteBuilder: knitr
Roxygen: list(markdown = TRUE)

Loading IDBacApp
Testing IDBacApp
✓ | OK F W S | Context
✓ | 1 |
✓ | 18 | 00_create_database [7.5 s]
✓ | 6 | 01_create_database [1.2 s]
✓ | 5 | test-createFuzzyVector
✓ | 1 | test-bootstrap [0.8 s]
✓ | 1 | brukerToMzml_popup
✓ | 1 | test-colored_dots [0.3 s]
✓ | 1 | test-colorpalette
✓ | 0 | controlBrukerDisplay

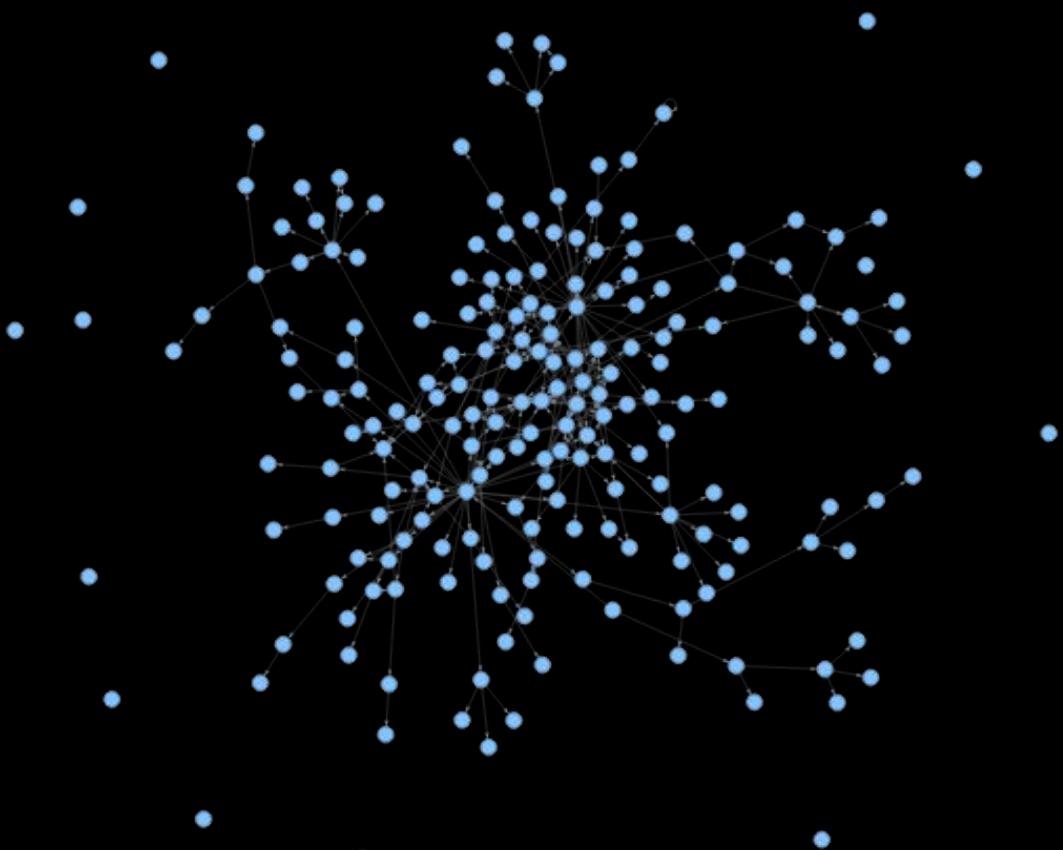
test-controlBrukerDisplay.R:18: skip: controlBrukerDisplay works
Reason: On Linux

✓ | 2 | copy_db
✓ | 1 | db_path_from_pool
✓ | 1 | distMatrix
✓ | 2 | findMSconvert
✓ | 6 | getMicrotyperFiles
✓ | 4 | getPeakData
✓ | 1 | test-hashr
✓ | 3 | test-insertLocale
✓ | 1 | test-labelsFromBrushedDendrogram
✓ | 9 | test-map384Well
✓ | 3 | mirrorplots
✓ | 1 | test-module_databaseTab
✓ | 1 | test-module_tsne
✓ | 4 | test-networkFromDF
✓ | 7 | parseDelimitedMS
✓ | 1 | test-pathSanitize
✓ | 1 | test-plotly_3d_scatter
✓ | 1 | test-runApp
✓ | 1 | test-sampleMapView
✓ | 3 | test-spectrumMatrixToMalDIQuant
✓ | 12 | sqlTableArchitecture
✓ | 3 | subtree
✓ | 6 | test-utils

== Results ==
Duration: 12.1 s

OK:    108
Failed:  0
Warnings: 1
Skipped: 1

```



Network of internal/external  
IDBac functions

Programming with IDBac

x +

127.0.0.1:4984

Programming With IDBac

1 Preamble

- 1.1 Major points:
- 1.2 Download IDBac example file

2 Connect to an IDBac Database

- 2.1 Connect to IDBac database

3 IDBac Databases Explained

- 3.1 Database Tables
- 3.2 How do I use this????

4 Add Data to an IDBac Database

- 4.1 Make a new empty IDBac database
- 4.2 Add data from an mzXML file
- 4.3 Adding multiple mzXML files

5 Starting With Bruker Data

6 Simple Analysis

- 6.1 Check database
- 6.2 What samples are in the database?

7 Working with IDBac from Python

IDBac Website

IDBac PNAS Publication

IDBac Video Protocol

Published with bookdown

Chapter 6 Simple Analysis | Prog

x +

127.0.0.1:4984/simple-analysis.html

Programming With IDBac

1 Preamble

- 1.1 Major points:
- 1.2 Download IDBac example file

2 Connect to an IDBac Database

- 2.1 Connect to IDBac database

3 IDBac Databases Explained

- 3.1 Database Tables
- 3.2 How do I use this????

4 Add Data to an IDBac Database

- 4.1 Make a new empty IDBac database
- 4.2 Add data from an mzXML file
- 4.3 Adding multiple mzXML files

5 Starting With Bruker Data

6 Simple Analysis

- 6.1 Check database
- 6.2 What samples are in the database?

7 Working with IDBac from Python

IDBac Website

IDBac PNAS Publication

IDBac Video Protocol

Published with bookdown

# Programming with IDBac

Chase Clark  
2020-02-15

## Chapter 1 Preamble

While some familiarity with R is suggested, the examples are enough for a novice to comfortably work through.

Suggestions or additions for content are welcome and can be submitted to [github.com/chasemc/programmingidbac](https://github.com/chasemc/programmingidbac). Note that this

### 1.1 Major points:

Things you want to do will revolve around:

- Creating an IDBac database from your raw (or converted) data
- Moving data from one IDBac database to another
- Accessing spectra
- Accessing peak-picked data
- Filtering data by some attribute

### 1.2 Download IDBac example

The data used in this book uses example data that can be found at <ftp://massive.ucsd.edu/MSV000084291>

```
library(here)
```

## Chapter 6 Simple Analysis

Load necessary packages for this tutorial:

If you haven't already, connect to an IDBac database as shown in "01\_connect-to-idbac-database".

Connect to the database

```
example_pool <- IDBacApp::idbac_connect(fileName = "idbac_experiment_file",
                                             filePath = here::here("data",
                                                       "example_data"))
```

```
my_plot <- IDBacApp::assembleMirrorPlots(sampleID1 = "172-7",
                                              sampleID2 = "172-10",
                                              peakPercentPresence = 0.7,
                                              lowerMassCutoff = 3000,
                                              upperMassCutoff = 15000,
                                              minSNR = 4,
                                              tolerance = 0.002,
                                              pool1 = example_pool$idbac_experiment_file,
                                              pool2 = example_pool$idbac_experiment_file)
```

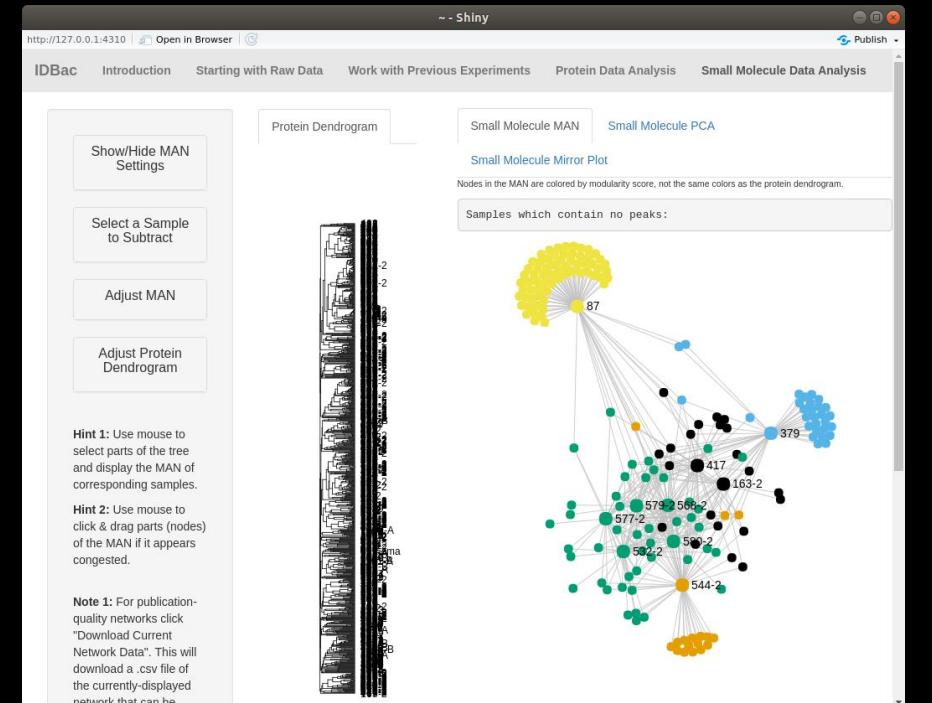
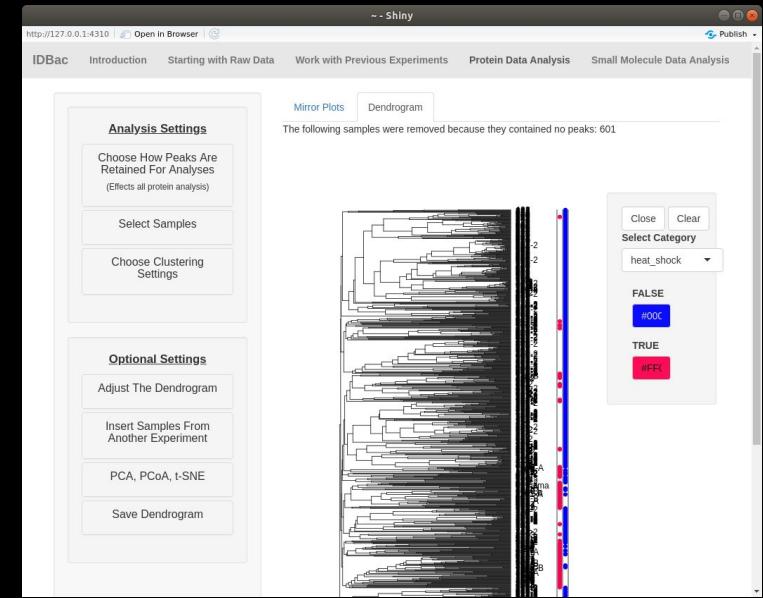
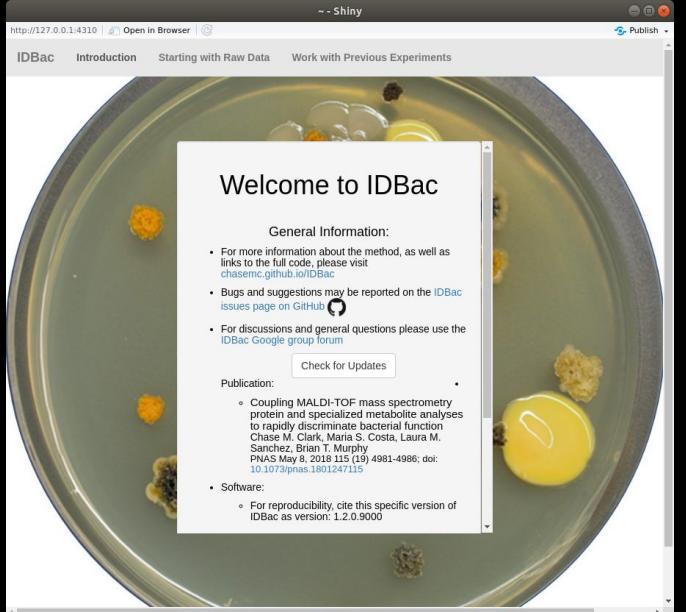
```
IDBacApp::mirrorPlot(my_plot)
```

Intensity

# IDBac

## MALDIquant mzR

## SQLite



# Coupling MALDI-TOF mass spectrometry protein and specialized metabolite analyses to rapidly discriminate bacterial function

Chase M. Clark<sup>1,2</sup>, Maria S. Costa<sup>1,2</sup>, Laura M. Sanchez<sup>1,2</sup>, and Brian T. Murphy<sup>1,2</sup><sup>1</sup>Department of Medicinal Chemistry and Pharmacognosy, College of Pharmacy, University of Illinois at Chicago, Chicago, IL; and <sup>2</sup>Faculty of Pharmaceutical Sciences, University of Iceland, Hagi, IS-107 Reykjavik, Iceland

Edited by Jeroil Meinwald, Cornell University, Ithaca, NY, and approved April 5, 2018 (received for review January 22, 2018)

For decades, researchers have used the ability to rapidly discriminate microbial communities with bacterial metabolism. Since specialized metabolites are critical to bacterial function and survival in the environment, we designed a data acquisition and bioinformatics technique (IDBac) that utilizes in situ matrix-assisted laser desorption ionization-mass spectrometry (MALDI-TOF MS) to analyze protein and specialized metabolite spectra recorded from single bacterial colonies picked from agar plates. We demonstrated the power of our approach by discriminating between two *Bacillus subtilis* strains in ~30 min using a workflow that included protein and metabolite extraction, peptide antibiotics surfactin and phasinatin, caused by a single frameshift mutation. Next, we used IDBac to detect subtle intra-species differences in the production of metal scavenging alycyclic polyketides of eight *Salinipora* spp. isolates. All isolates had isolates that share >95% sequence similarity in the 16S rRNA gene.Finally, we used IDBac to simultaneously extract protein and specialized metabolite MS profiles from unidentified Lake Michigan sponge-associated bacteria isolated from an agar plate. In just 3 h, we created hierarchical protein MS groupings of 11 environmental isolates that were highly correlated with their species and accurately mirrored phylogenetic groupings. We further distinguished isolates within these groupings, which share nearly identical 16S rRNA gene sequence identity, based on interspecies and intraspecies differences in their specialized metabolite production. IDBac is an attempt to couple *in situ* MS analysis of protein content and specialized metabolite production to allow for facile discrimination of closely related bacterial isolates.

## Significance

Mass spectrometry is a powerful technique that has been used to identify bacteria by their protein content and to assess potential metabolic diversity by the analysis of their specialized metabolites. However, until now these analyses have operated independently, which has resulted in the inability to rapidly connect bacterial phylogenetic identity with potential environmental function. To bridge this gap, we designed a MALDI-TOF mass spectrometry-based pipeline (IDBac) to integrate data from both intact protein and specialized metabolite spectra directly from bacterial cells grown on agar. This technique organizes bacteria into highly similar phylogenetic groups and allows for comparison of metabolic differences of hundreds of isolates in just a few hours.

<https://doi.org/10.1073/pnas.1801247115>

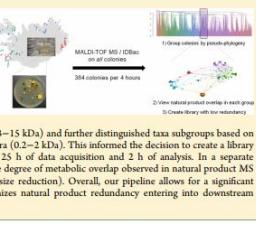
**JOURNAL OF NATURAL PRODUCTS** Article Cite This: J. Nat. Prod. XXXX, XXX, XXX–XXX [pubs.acs.org/jnp](http://pubs.acs.org/jnp)

**Minimizing Taxonomic and Natural Product Redundancy in Microbial Libraries Using MALDI-TOF MS and the Bioinformatics Pipeline IDBac**

Maria S. Costa,<sup>1,2,3</sup> Chase M. Clark,<sup>2,4</sup> Sessela Ómarsdóttir,<sup>2</sup> Laura M. Sanchez,<sup>2,5</sup> and Brian T. Murphy<sup>2,6,7</sup>

<sup>1</sup>Faculty of Pharmaceutical Sciences, University of Iceland, Hagi, Hofvallagata 53, IS-107 Reykjavík, Iceland  
<sup>2</sup>Department of Pharmaceutical Sciences, College of Pharmacy, University of Illinois at Chicago, 833 South Wood Street (MC 781), Room 539, Chicago, Illinois 60607, United States  
<sup>3</sup>Supporting Information

**ABSTRACT:** Libraries of microorganisms have been a cornerstone of drug discovery efforts since the mid-1950s, but strain duplication in some libraries has resulted in unwanted natural product redundancy. In the current study, we implemented a workflow that minimizes both the natural product overlap and the total number of bacterial isolates in a library. Using a collection of 86 environmental samples, we purified every distinct bacterial colony off solid media plates derived from 86 environmental samples. We employed our mass spectrometry (MS)-based IDBac workflow on these isolates to form groups of taxa based on protein MS fingerprints (3–15 kDa) and further distinguished taxa subgroups based on their degree of overlap within corresponding natural product spectra (0.2–2 kDa). This informed the decision to create a library of 301 isolates spanning 54 genera. This process required over 25 h of manual labor and 2 of automation. In a separate experiment, we reduced the size of the library by 16% to 250 isolates by purifying every distinct bacterial colony (from 86 environmental samples) (from 833 to 233 isolates, a 72.0% size reduction). Overall, our pipeline allows for a significant reduction in costs associated with library generation and minimizes natural product redundancy entering into downstream biological screening efforts.



<https://doi.org/10.1021/acs.jnatprod.9b00168>

# IDBac Publications

<https://chasemc.github.io/IDBac>



EDITORIAL



## A Call to Action: the Need for Standardization in Developing Open-Source Mass Spectrometry-Based Methods for Microbial Subspecies Discrimination

Chase M. Clark,<sup>a</sup> Brian T. Murphy,<sup>a</sup> Laura M. Sanchez<sup>b</sup><sup>a</sup>Department of Pharmaceutical Sciences, University of Illinois at Chicago, Chicago, Illinois, USA**KEYWORDS** MALDI-TOF MS, bioinformatics, dereplication, microbial ecology

In the last decade, there has been a renewed push by academic researchers to create rapid and accurate techniques to differentiate, identify, and prioritize culturable microbial isolates. One such technique that continues to gain momentum among microbiologists is matrix-assisted laser desorption-ionization time of flight mass spectrometry (MALDI-TOF MS). It is an established, inexpensive technique commonly used to rapidly identify microbial taxa and differentiate culturable microbes. This technology has become commonplace in clinical and veterinary laboratories where rigorously validated methods are used in conjunction with commercially available reference databases to identify pathogenic microorganisms. However, the broader community, especially laboratories working with environmental microbes, typically cannot access the expensive software and databases. It is our opinion that this community, which relies on free and open-source software, currently lacks a coherent set of accepted experimental practices, including employment of internal standard strains, statistically driven determination of biological and technical replicates, and deposition of MS data into open-access repositories. Establishing guidelines would enable researchers to better compare microbial typing methods and advance our ability to group and delineate environmental isolates in an effective manner, particularly at the subspecies level.

<https://doi.org/10.1128/mSystems.00813-19>

**jove** Search 10,747 video articles Advanced

BIOPHARMA NEW ABOUT Jove FOR LIBRARIANS PUBLISH VIDEO JOURNAL SCIENCE EDUCATION

ABSTRACT INTRODUCTION PROTOCOL RESULTS DISCUSSION MATERIALS REFERENCES DOWNLOADS

**BIOCHEMISTRY**

Using the Open-Source MALDI TOF-MS IDBac Pipeline for Analysis of Microbial Protein and Specialized Metabolite Data

Chase M. Clark<sup>1</sup>, Maria S. Costa<sup>1,2</sup>, Erin Conley<sup>1</sup>, Emma Li<sup>1</sup>, Laura M. Sanchez<sup>1</sup>, Brian T. Murphy<sup>1</sup>

<sup>1</sup>Department of Medicinal Chemistry and Pharmacognosy, College of Pharmacy, University of Illinois at Chicago, <sup>2</sup>Faculty of Pharmaceutical Sciences, University of Iceland

This content is open access.

CHAPTERS

0.04 Title  
 1.10 Preparation of MALDI Target Plates and Data Acquisition  
 2.02 Installing the IDBac Software and Starting with Raw Data  
 2.47 Work with Previous Experiments  
 3.22 Setting up Protein Data Analysis and Creating Mirror Plots  
 4.07 Clustering Samples Using Protein Data  
 4.49 Customizing the Protein Dendrogram and Inserting Samples from a Separate Experiment into the Dendrogram  
 5.32 Analyzing Specialized Metabolite Data and Metabolite Association Networks (MANs)

RESULTS

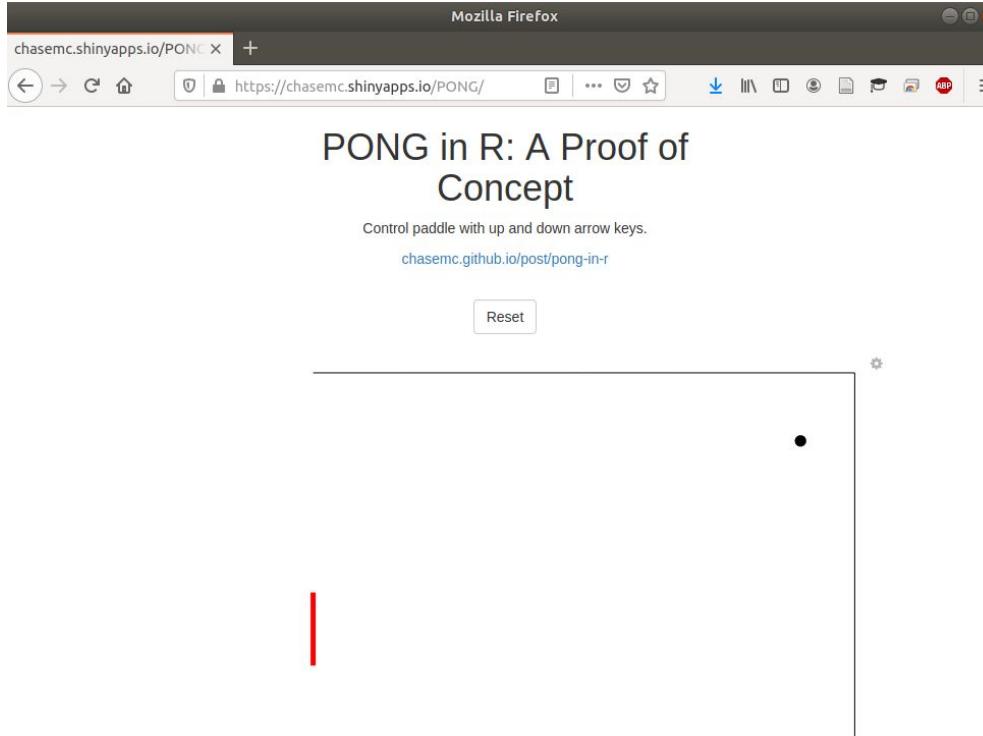
ISSUE 147 DOI: 10.3791/59219 PUBLISHED: 5/15/2019 0 COMMENTS PDF EMBED ADD TO FAVORITES

<https://doi.org/10.3791/59219>

<https://doi.org/10.3791/59219>

# How will users access your software?

## Web Deployment (Server)



<https://chasemc.shinyapps.io/PONG/>

## Local GUI

A screenshot of a GitHub repository page for "mzfinder". The repository has 13 files and 1 directory listed, all updated "4 months ago". The files include .gitignore, CODE\_OF\_CONDUCT.md, DESCRIPTION, Dockerfile, LICENSE, LICENSE.md, NAMESPACE, NEWS.md, README.md, app.R, npsearch.Rproj, and README.md. The README.md file contains installation instructions:

```
install.packages("remotes") # Run this if you don't have the remotes package
remotes::install_github("chasemc/mzfinder")

# And, to run the app:
mzfinder::run_app()
```

A note at the bottom states: "Please note that the 'mzfinder' project is released with a [Contributor Code of Conduct](#). By contributing to this project, you agree to abide by its terms."

<https://github.com/chasemc/mzfinder>

# How will users access your software?

## Web Deployment

large files + long computations



# How will users access your software?

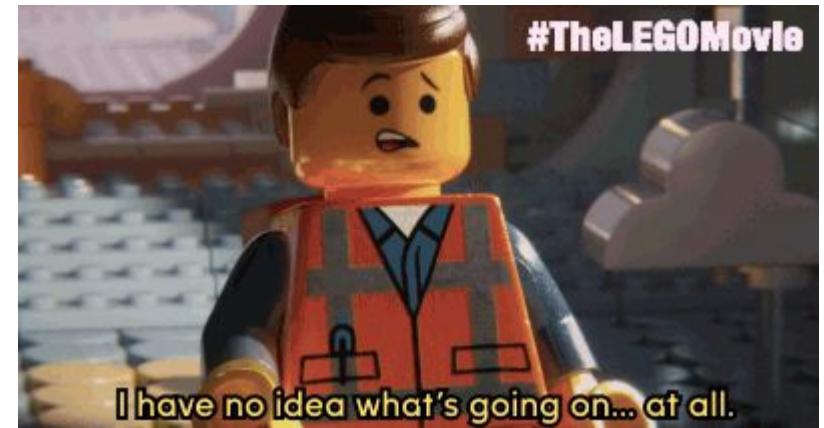
## Web Deployment

large files + long computations



## Local Installation

Large percent of users won't even attempt  
Local technical support





# {electricShine}

## Installable Shiny/Electron Apps

• lifecycle experimental

Windows CI:

• build passing

Mac and Linux CI

• build passing

Is easy! One meta-function:

```
electricShine::buildElectricApp(  
  app_name = "My_App",  
  description = "My demo application",  
  package_name = "demoApp",  
  semantic_version = "1.0.0",  
  build_path = buildPath,  
  mran_date = MRANdate,  
  function_name = "run_app",  
  github_repo = "chasemc/demoApp",  
  local_path = NULL  
)
```

Self-Sufficient

- R- from specified MRAN date
  - R packages- from specified MRAN date
  - Your Shiny app!
- 
- Currently: Windows/Mac
  - Future: Linux
  - Compatible with continuous-deployment

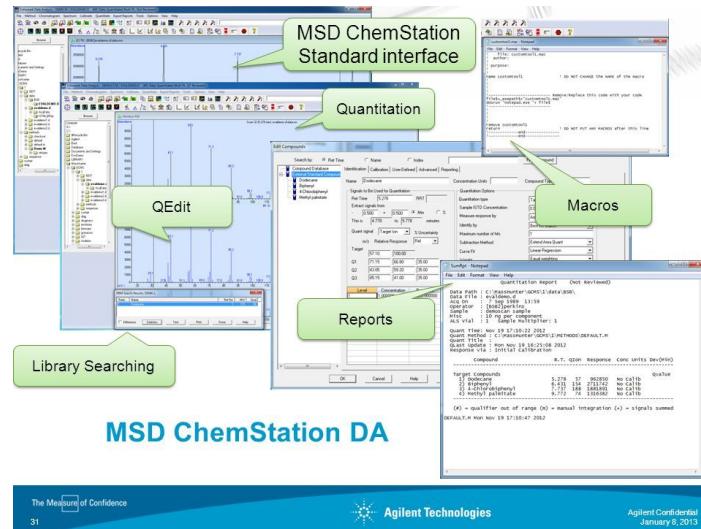
<https://github.com/chasemc/electricShine>



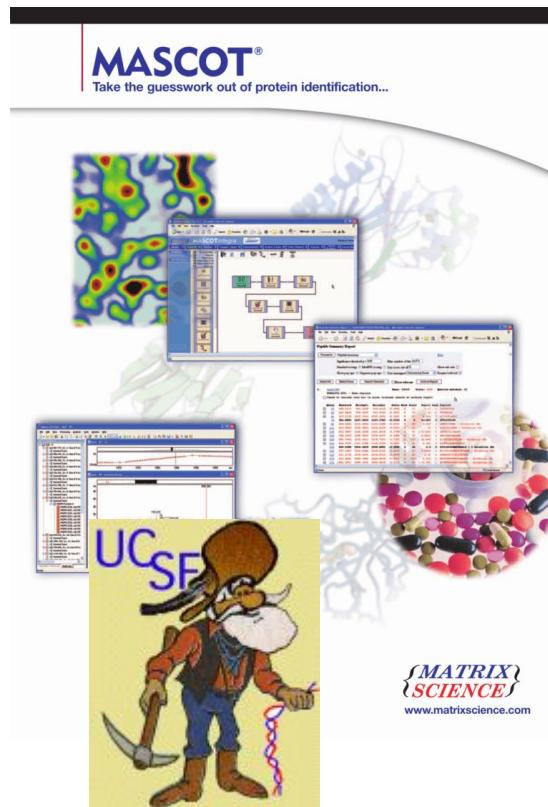
@ChasingMicrobes

# The many faces of mass spectrometry software

## GUI



## Server



## Command line

This screenshot shows a command line interface (cmd.exe) running on Windows. The title bar says "C:\WINDOWS\System32\cmd.exe". The window displays the usage instructions for the CompassXport command line utility version 1.2. The text is as follows:

```
C:\Documents and Settings\AndrewG>compassxport
CompassXport command line utility version 1.2
command line usage

Required parameters:
1) Convert a single raw data file into a single mzXML file.
If no output file name is specified the output is placed
next to the raw data file
[-a input file name] [-o output file name]

2) Convert multiple raw data files into multiple mzXML files.
The output files are placed next to the raw data files.
File name of output files can be specified (analysis.mzXML is default)
[-multi root path of analyses] [-multiName file name of output files]

Optional Parameters:
3) Set log level: [-log level]
where level can be: <non, error, all>
Default: all

4) Switch to control raw data export: [-raw X]
X can be 1 <export as raw data> or 0 <convert raw data to line spectra>
Default: 0

Input file name can be:
- analysis.baf acquisition data files
- analysis.yep acquisition data files
- AutoExecute run file for LCMaldi
- fid acquisition data files
```

# The many faces of mass spectrometry software

GUI



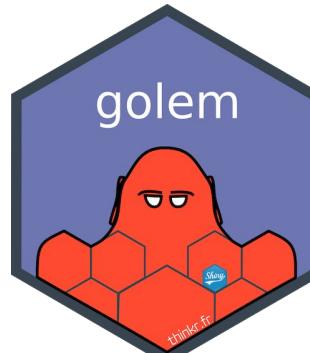
Server



Command line



{plumber} converts your existing R code to a web API using a handful of special one-line comments.



{golem} is an opinionated framework for building production-grade shiny applications.

# SCIENCE TODAY

ICELAND -  
THE FUTURE  
OF MEDICINE?

SCIENCE  
MUSEUM

## Funding:

National Center For  
Complementary & Integrative  
Health of the National Institutes  
of Health Award Number:  
F31AT010419



## Murphy Lab Current Grad and Post-Doc Members:

Dr. Brian Murphy  
Chase Clark  
Antonio Hernandez  
Dr. Jeongho Lee  
Dr. Tuan Anh Tran  
Dr. Linh Nguyen  
Maryam Elfeki



@ChasingMicrobes  
chasemc.github.io

MEDICINAL  
CHEMISTRY  
AND  
PHARMACOGNOSY  
COLLEGE  
OF PHARMACY

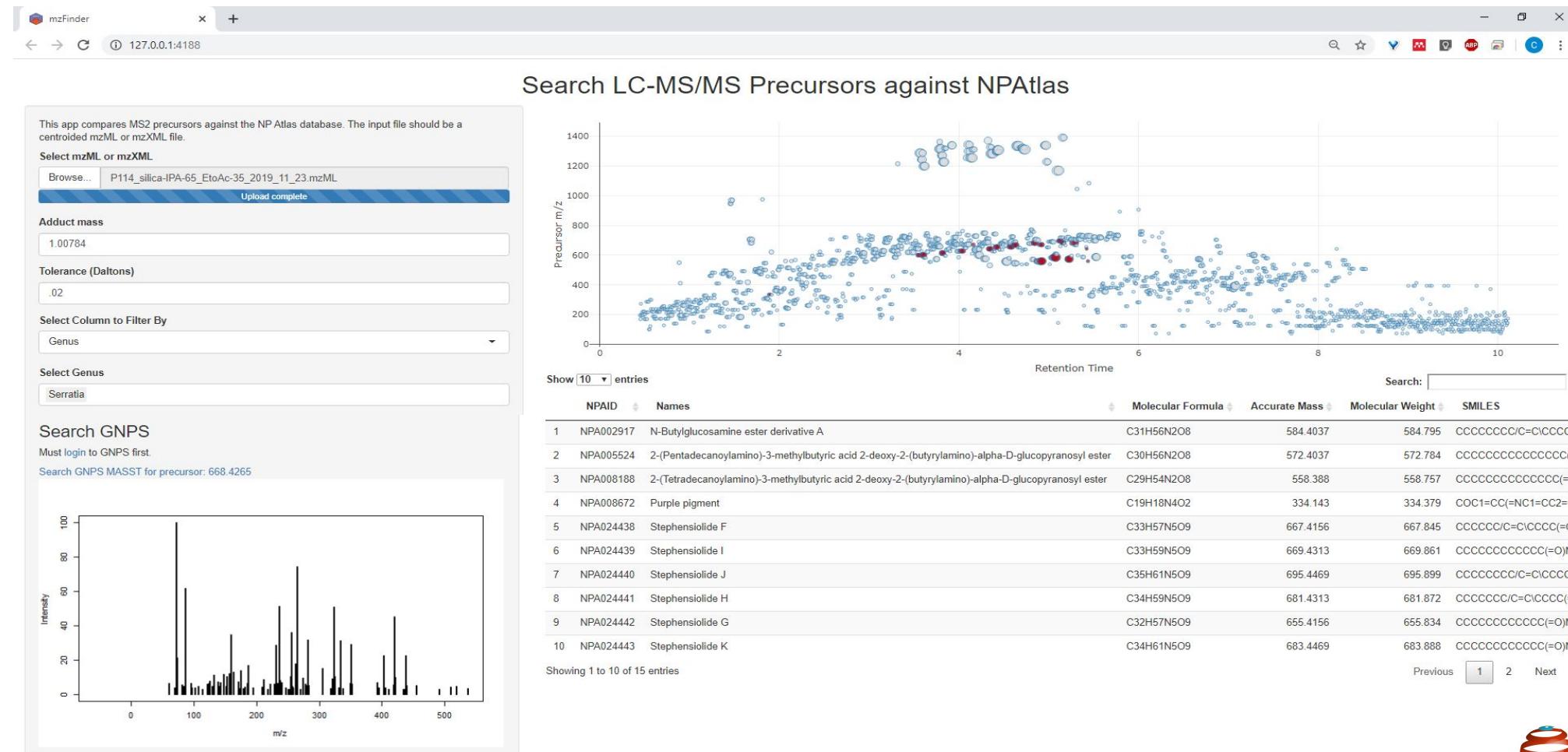




@ChasingMicrobes  
chasemc.github.io

[github.com/chasemc/presentations/may\\_institute/2020](https://github.com/chasemc/presentations/may_institute/2020)

# Is R Fast Enough? Usually Yes



<https://github.com/chasemc/mzfinder>

