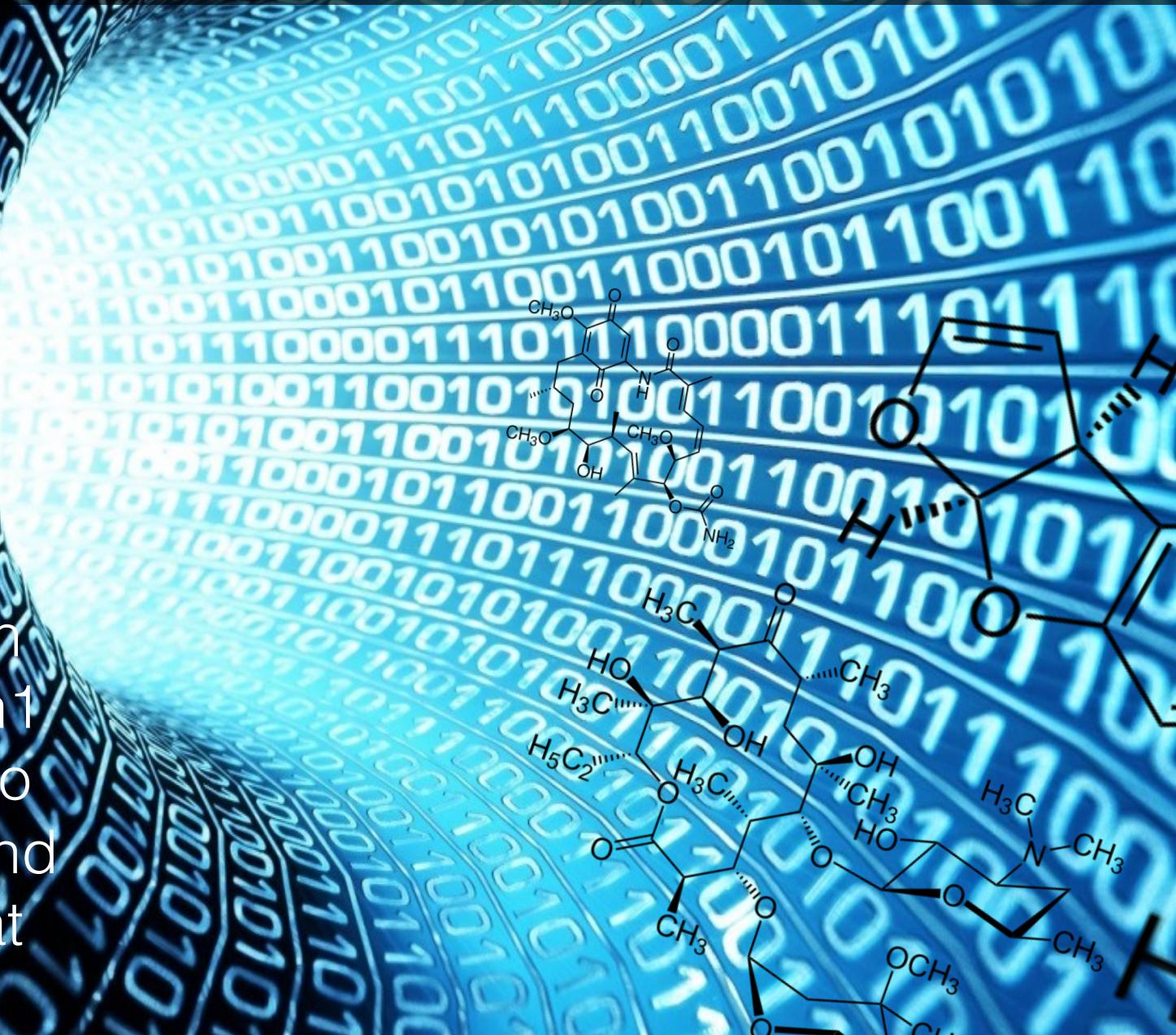
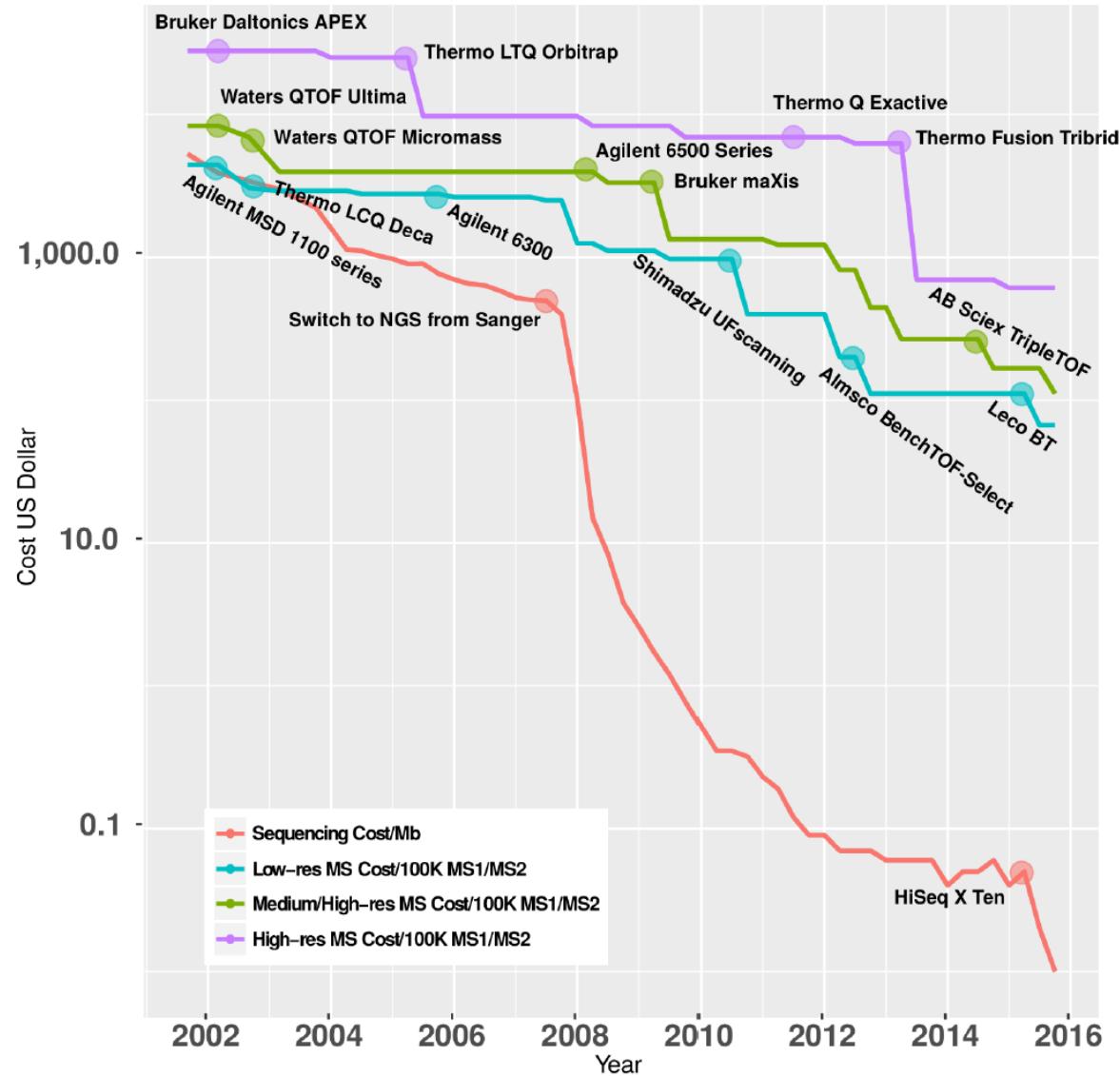


# The past, present and future of data in small molecule analysis, including metabolomics, and what impact standards will have.

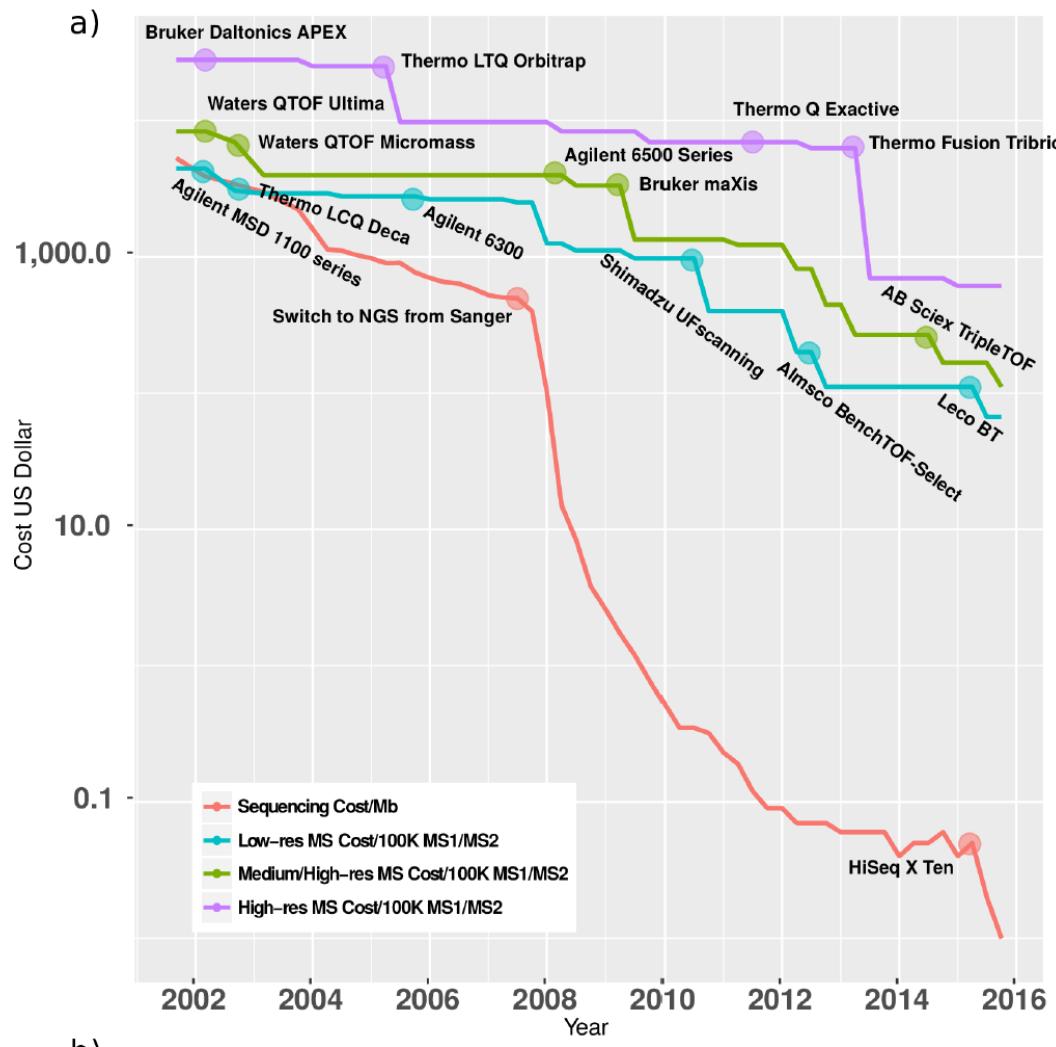
Pieter C Dorrestein  
Twitter: Pdorrestein1  
You are welcome to  
take screenshots and  
share the work that  
is presented



# We are entering an amazing time in mass spectrometry

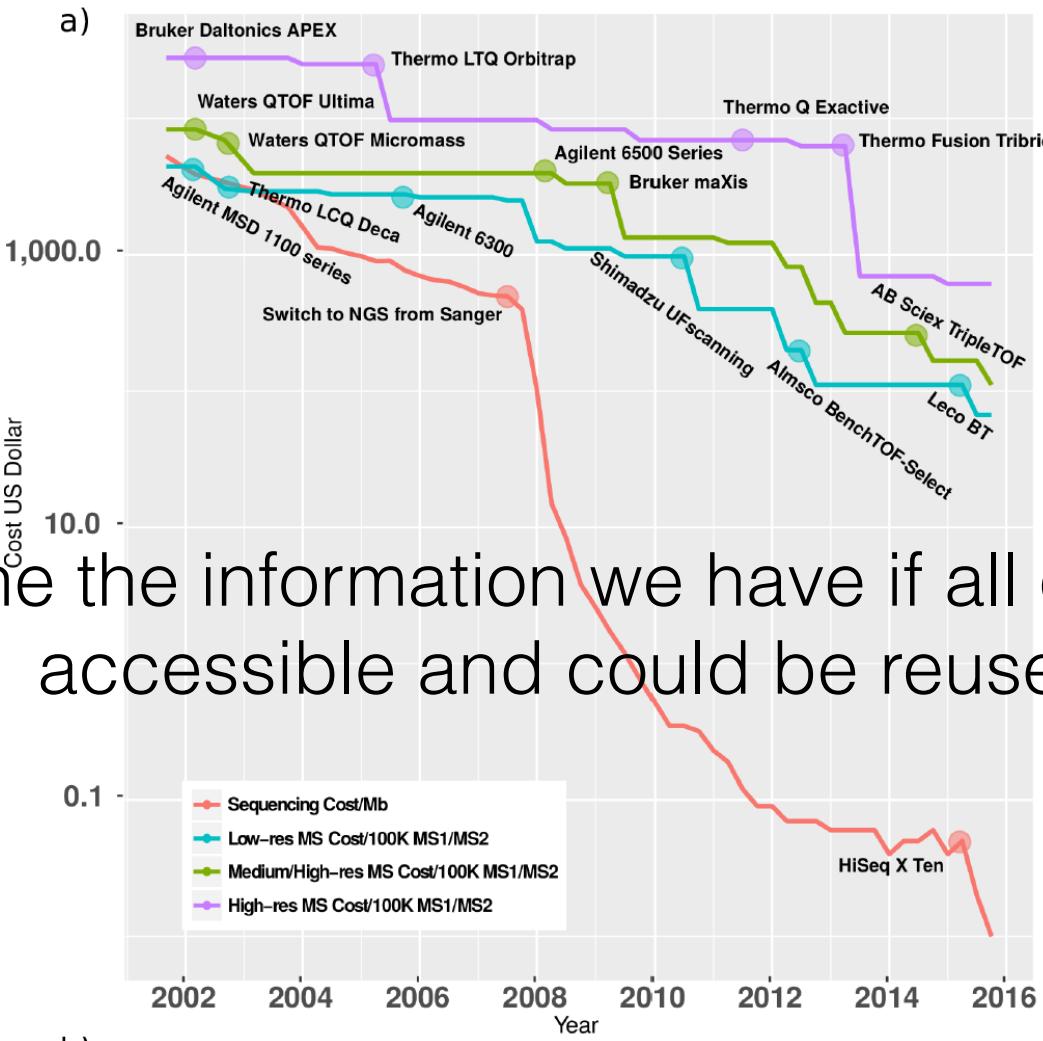


# We are entering an amazing time in mass spectrometry



We can collect more high quality, high resolution data than ever before

# We are entering an amazing time in mass spectrometry



Just imagine the information we have if all of the data was accessible and could be reused?

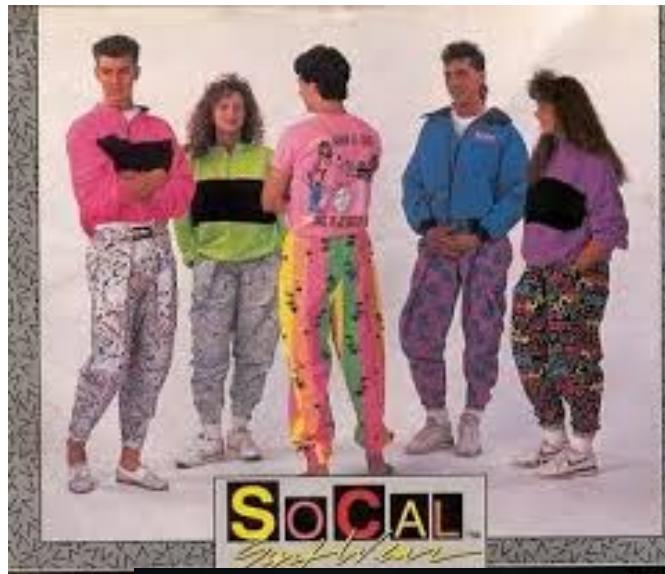
We can collect more high quality, high resolution data than ever before

There are ~ a dozen papers in the literature that reuse small molecule mass spectrometry data sets, and ~50% of them are for my own lab.

Why??

Different communities have different standards and are far behind compared to sequencing.

Comparatively IMO still stuck in...



# What if we ask this of molecules? where do we go?

The image shows the NCBI BLAST search interface on the left and a conceptual diagram of a retrieval infrastructure on the right.

**NCBI BLAST Search Interface:**

- Header:** NIH U.S. National Library of Medicine, NCBI National Center for Biotechnology Information, Sign in to NCBI
- Title:** BLAST® » blastp suite
- Navigation:** Home, Help, Reset page, Bookmark
- Search Options:** Enter Query Sequence, Enter accession number(s), gi(s), or FASTA sequence(s), Or, upload file, Choose File (No file chosen), Job Title, Enter a descriptive title for your BLAST search, Align two or more sequences
- Choose Search Set:** Database (Non-redundant protein sequences (nr)), Organism (Optional), Exclude (Models (XM/XP), Uncultured/environmental sample sequences), Entrez Query (Optional)
- Sequence View:** A large sequence snippet is shown: GACAGACATGACTTGGATTCCCCAGGAGGAGTTGGCAACCCATTCCAAGAGCTTGAACCCCTCTGAGAAATTCATCTCAGCACA... (truncated)

**Conceptual Diagram (right):**

- Database:** A large cylinder labeled "Database".
- Knowledgebase:** A smaller cylinder labeled "Knowledgebase".
- Magnifying Glass:** A magnifying glass is positioned over the Database cylinder, indicating the search process.
- Text Overlay:** Where is this seq. found? any analogs? + possible functions
- YouTube Link:** YouTube Create custom database
- Entrez Query:** Enter an Entrez query to limit search

## Retrieval infrastructure

What if we ask this of molecules? where do we go?



Different communities have different standards and are far behind compared to sequencing.

Comparatively IMO still stuck in...

### **Basic Local Alignment Search Tool**

Stephen F. Altschul<sup>1</sup>, Warren Gish<sup>1</sup>, Webb Miller<sup>2</sup>  
Eugene W. Myers<sup>3</sup> and David J. Lipman<sup>1</sup>

<sup>1</sup>*National Center for Biotechnology Information  
National Library of Medicine, National Institutes of Health  
Bethesda, MD 20894, U.S.A.*

<sup>2</sup>*Department of Computer Science  
The Pennsylvania State University, University Park, PA 16802, U.S.A.*

<sup>3</sup>*Department of Computer Science  
University of Arizona, Tucson, AZ 85721, U.S.A.*

*(Received 26 February 1990; accepted 15 May 1990)*

A new approach to rapid sequence comparison, basic local alignment search tool (BLAST), directly approximates alignments that optimize a measure of local similarity, the maximal segment pair (MSP) score. Recent mathematical results on the stochastic properties of MSP scores allow an analysis of the performance of this method as well as the statistical significance of alignments it generates. The basic algorithm is simple and robust; it can be implemented in a number of ways and applied in a variety of contexts including straightforward DNA and protein sequence database searches, motif searches, gene identification searches, and in the analysis of multiple regions of similarity in long DNA sequences. In addition to its flexibility and tractability to mathematical analysis, BLAST is an order of magnitude faster than existing sequence comparison tools of comparable sensitivity.

..and there were 1000s algorithms developed (and are still being developed.)

There are a variety of endeavors (not an exhaustive list below) that are aiming to tackle data analysis and/or the data sharing and analysis reproducibility problem.



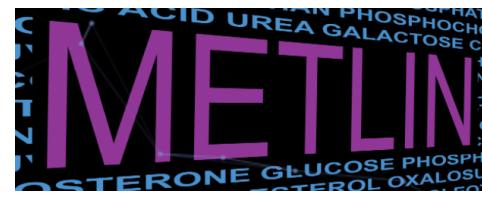
GNPS



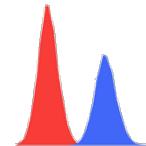
Golm Metabolome Database



MetaboLights



MoNA - MassBank of North America



MZmine 2



PhenoMeNal ReSpect

Large-Scale Computing for Medical Metabolomics



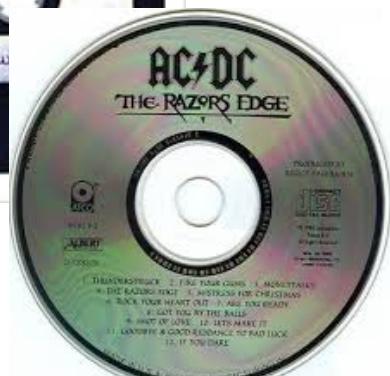
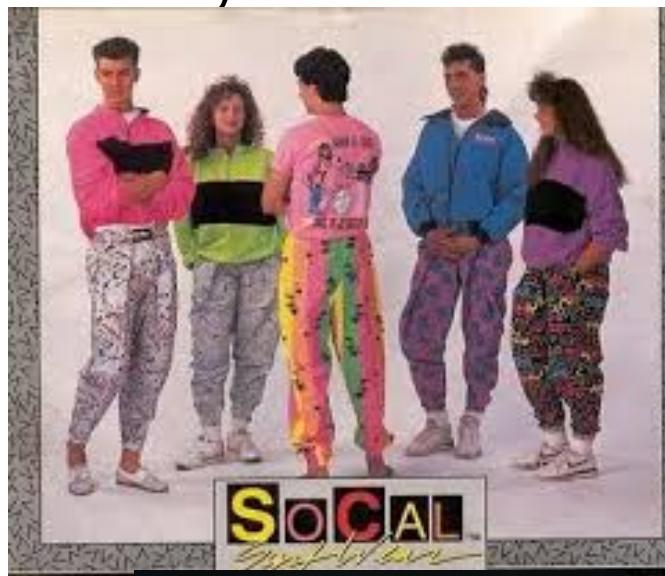
OpenMSI



No expectation for data and or knowledge sharing in digital format.

Of ~20 labs I have asked to share data, 0 have been able to.

Comparatively IMO still stuck in...



[Browse the Journal](#)[Articles ASAP](#)[Current Issue](#)[Submission & Review](#)[Open Access](#)[About the Journal](#)

## Article

 [Previous](#)

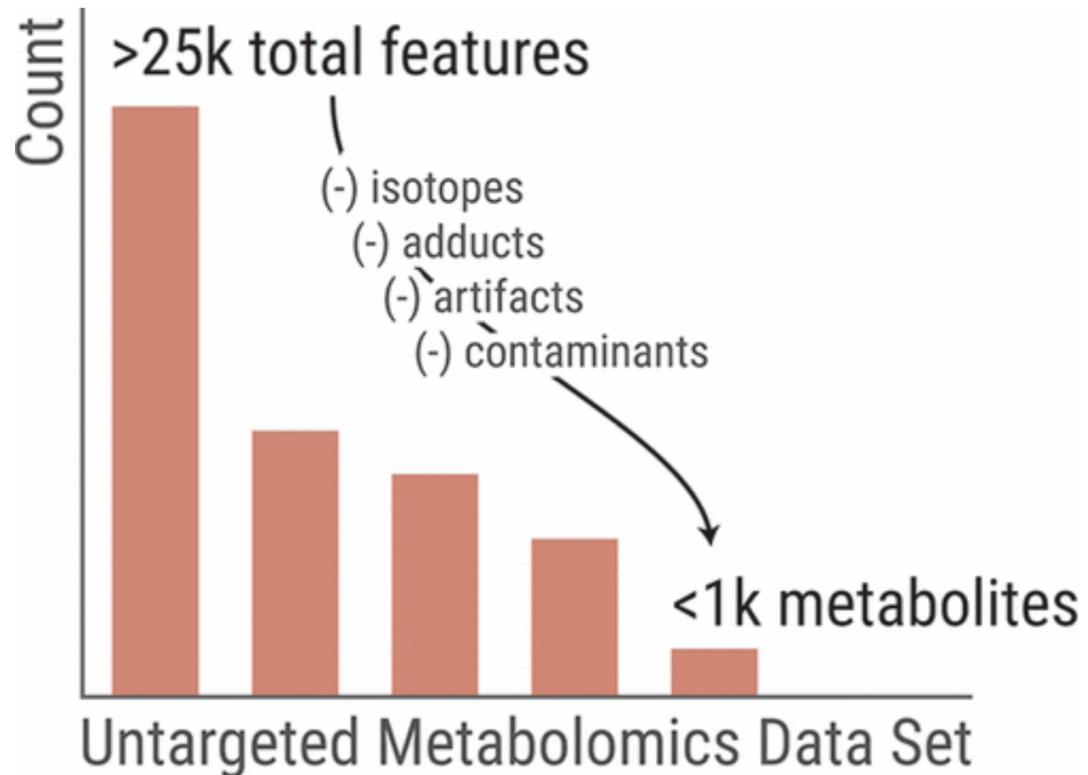
# Systems-Level Annotation of a Metabolomics Data Set Reduces 25000 Features to Fewer than 1000 Unique Metabolites

*Anal. Chem.*, 2017, 89 (19), pp 10397–10406

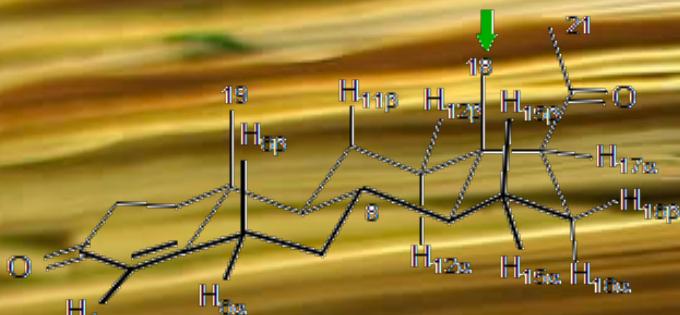
DOI: 10.1021/acs.analchem.7b02380

Publication Date (Web): September 15, 2017

Copyright © 2017 American Chemical Society



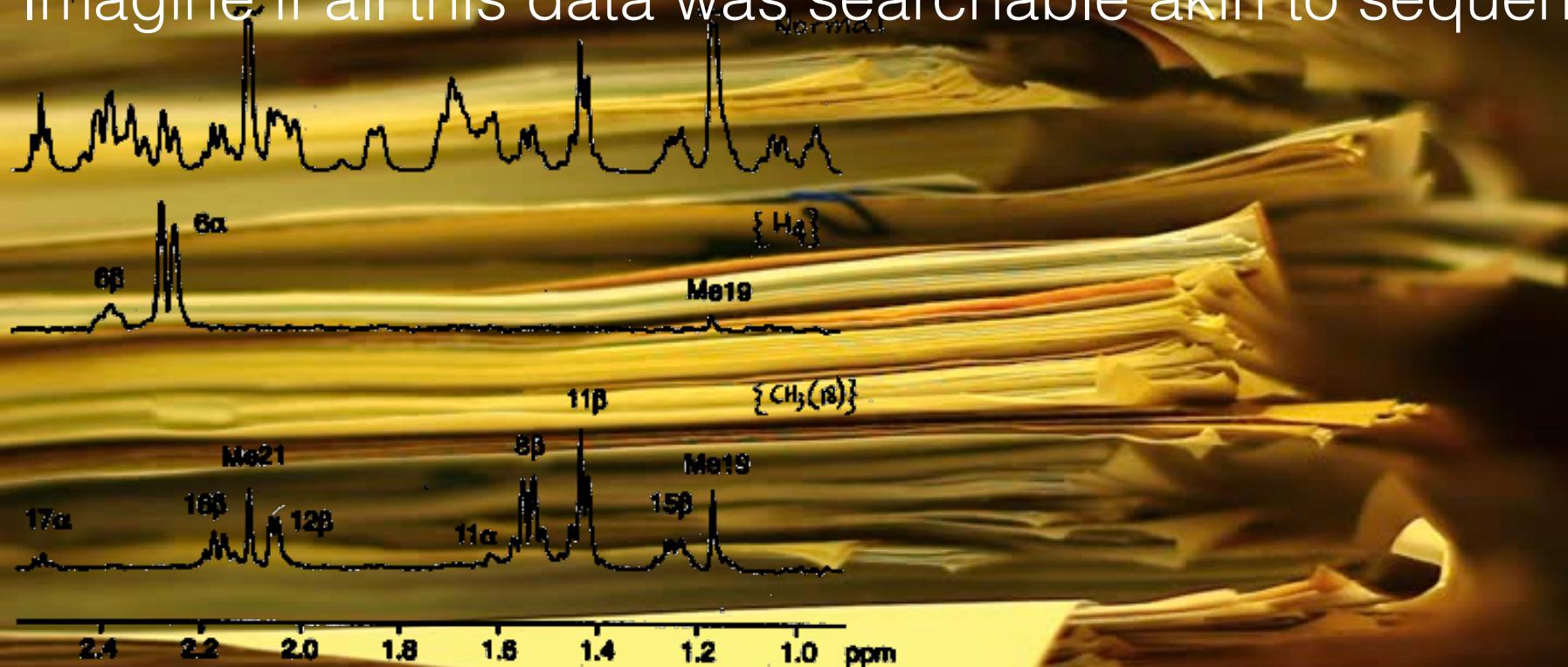
# What do you get when you ask a natural product scientist on data from their extract or molecule?



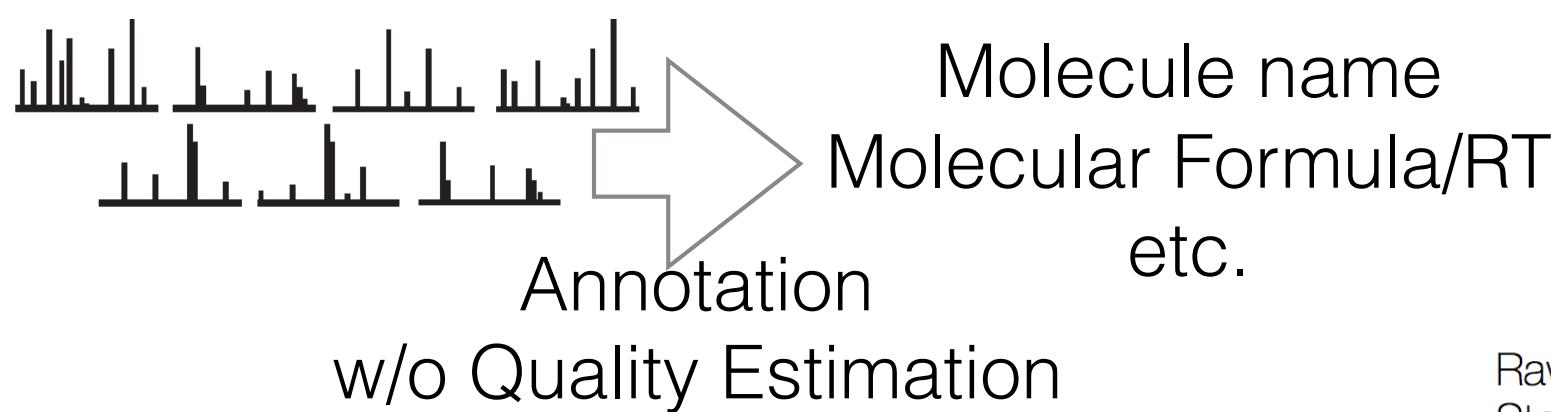
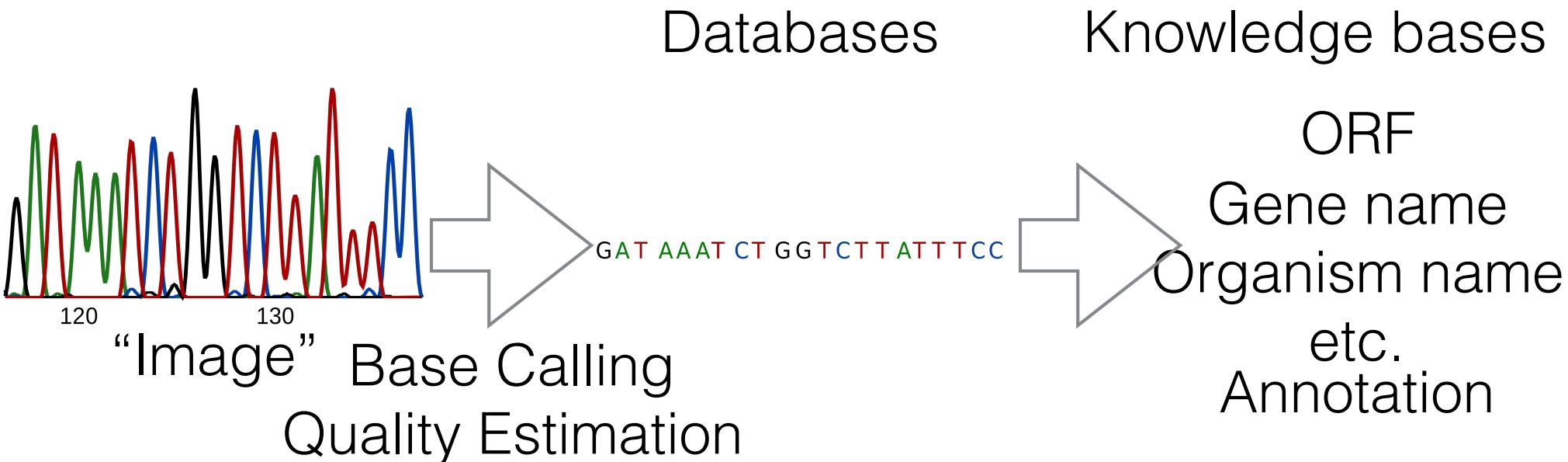
Where does it end up?

- Tables and/or images in Papers
- And Supporting information figures
- Structure in chemical databases

Imagine if all this data was searchable akin to sequence data?



I think raw data deposition is critical-not everyone agrees in the MS community. People that disagree rationalize what is done in sequencing.

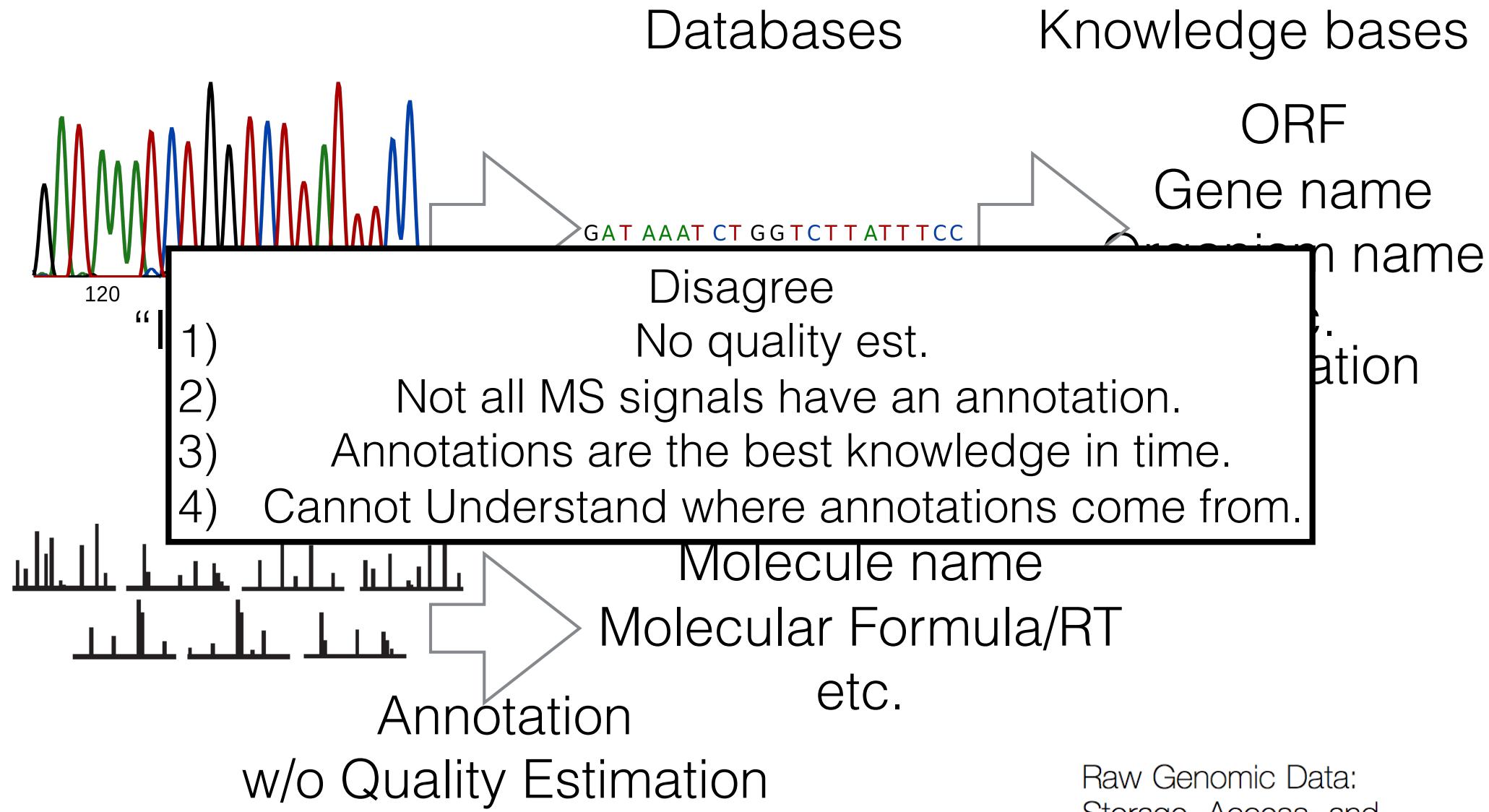


## On the Future of Genomic Data

Raw Genomic Data:  
Storage, Access, and  
Sharing

Mahsa Shabani,<sup>1,\*</sup>  
Danya Vears,<sup>1</sup> and  
Pascal Borry<sup>1</sup> Cell press reviews 2018

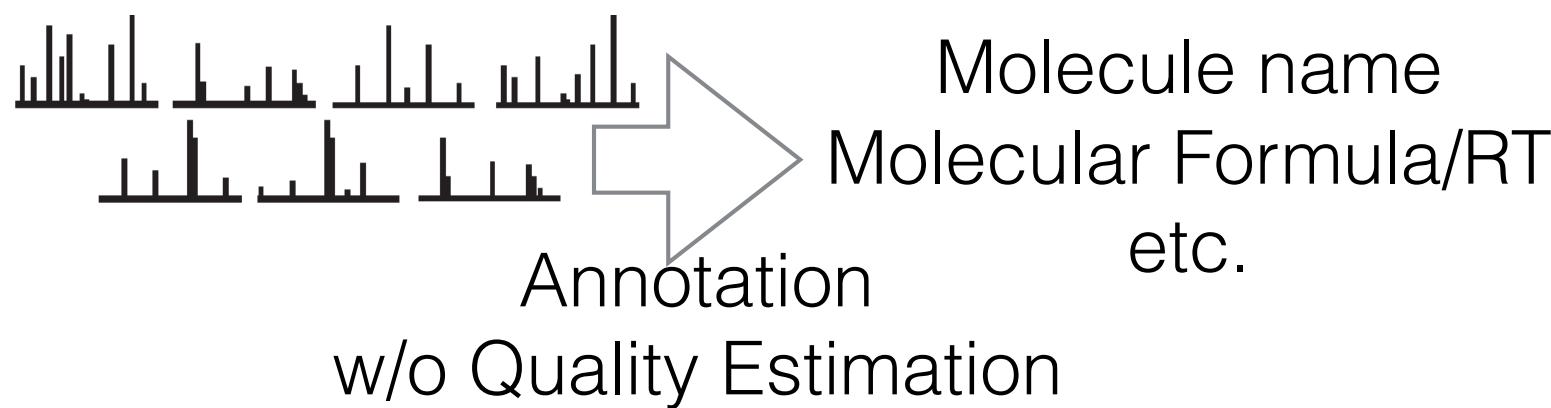
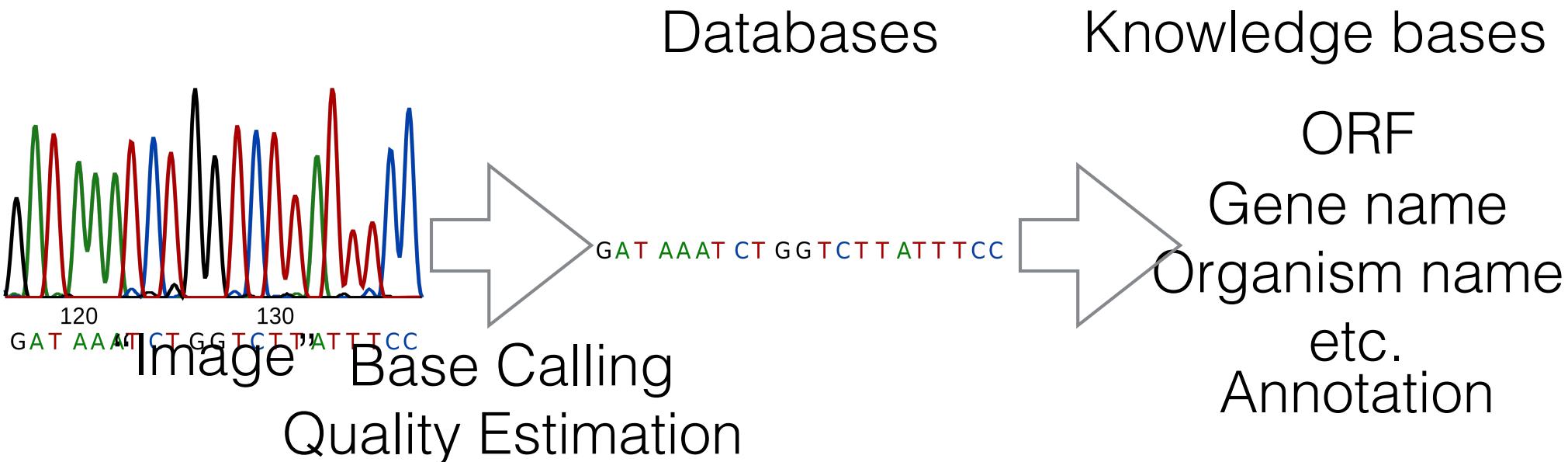
I think raw data deposition is critical-not everyone agrees in the MS community. People that disagree rationalize what is done in sequencing.



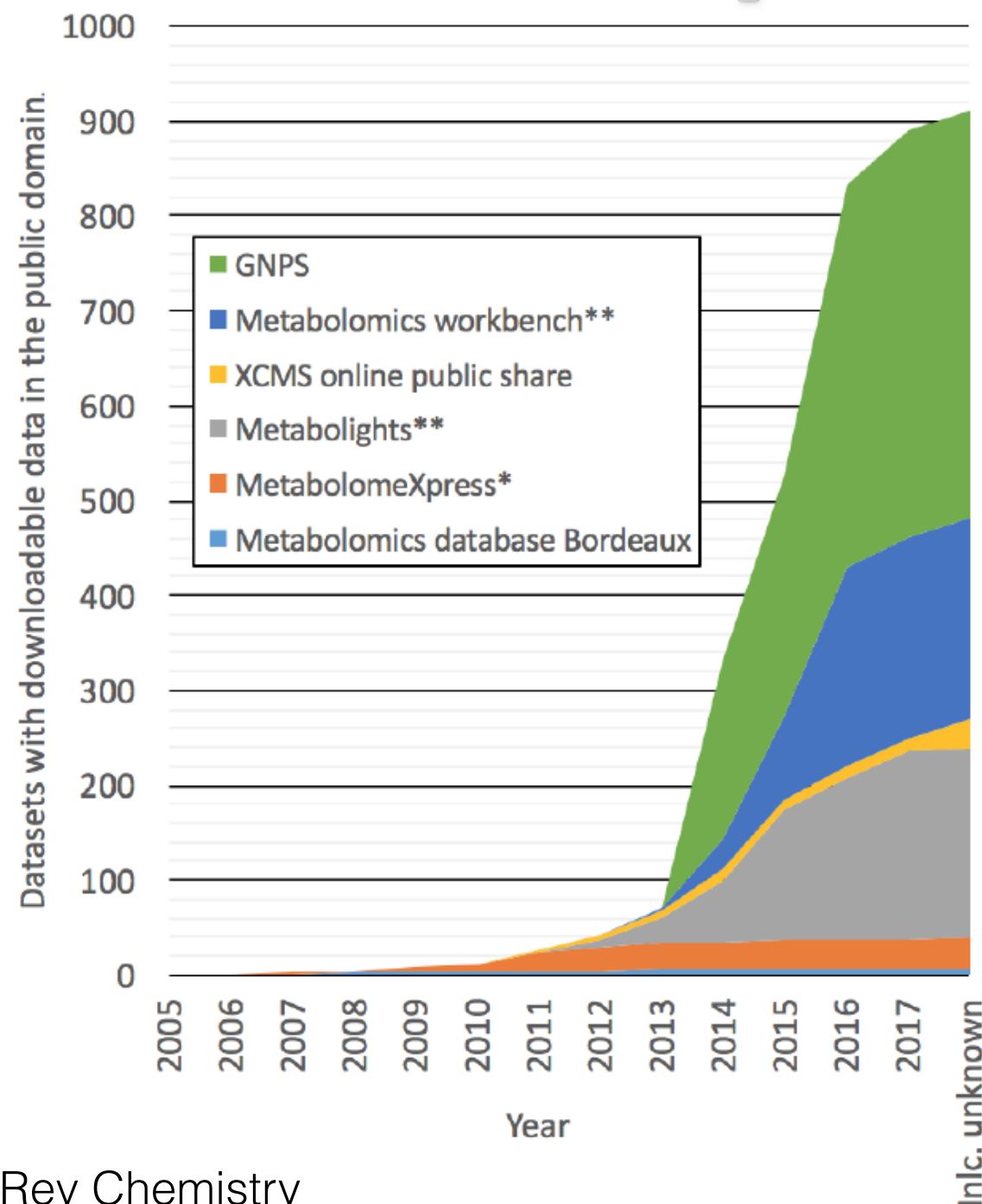
Raw Genomic Data:  
Storage, Access, and  
Sharing

Mahsa Shabani,<sup>1,\*</sup>  
Danya Vears,<sup>1</sup> and  
Pascal Borny<sup>1</sup> Cell press reviews 2018

I think raw data deposition is critical-not everyone agrees in the MS community. People that disagree rationalize what is done in sequencing.



# Metabolomics Projects with Downloadable Data in the Public Domain March 7, 2017

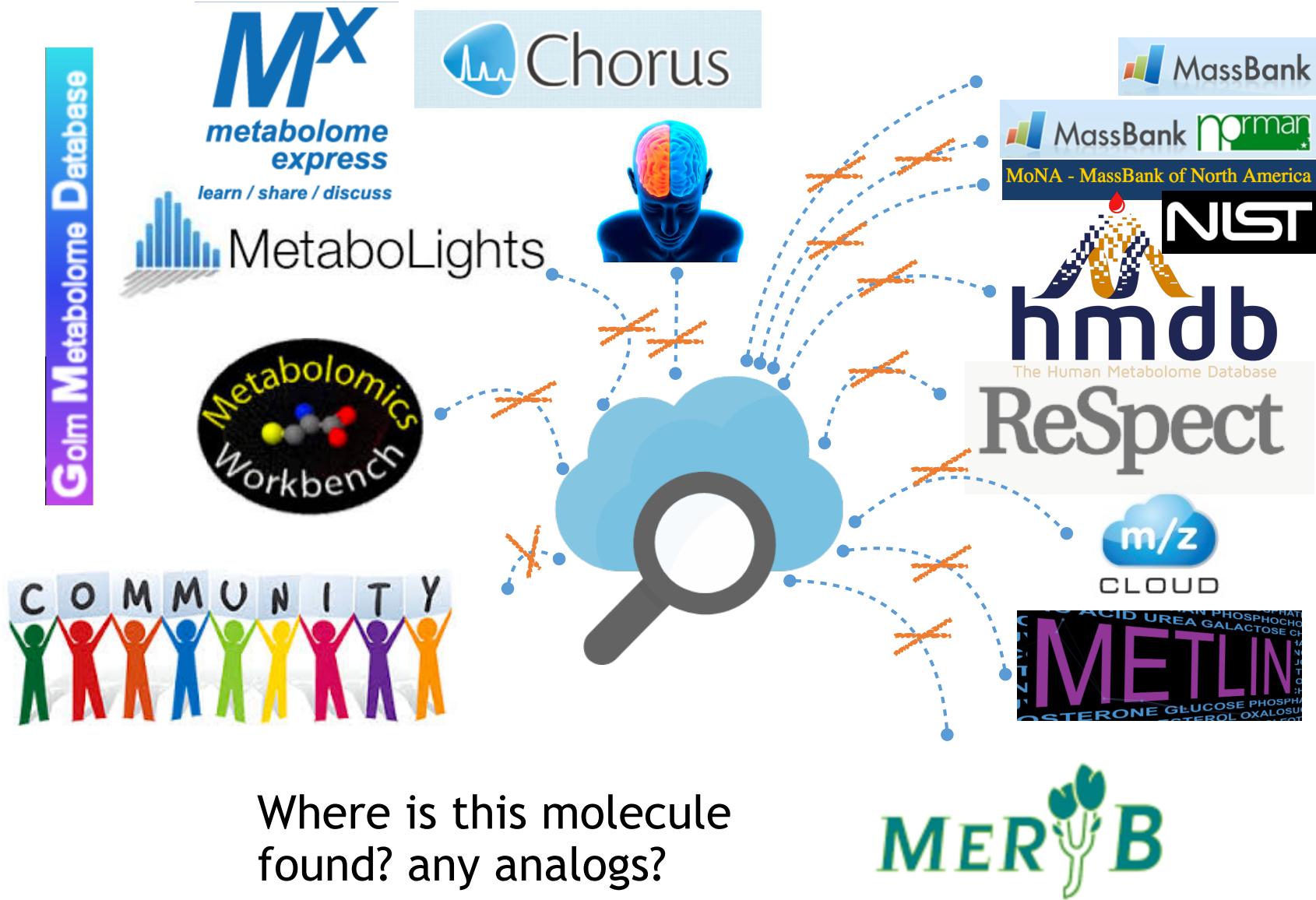




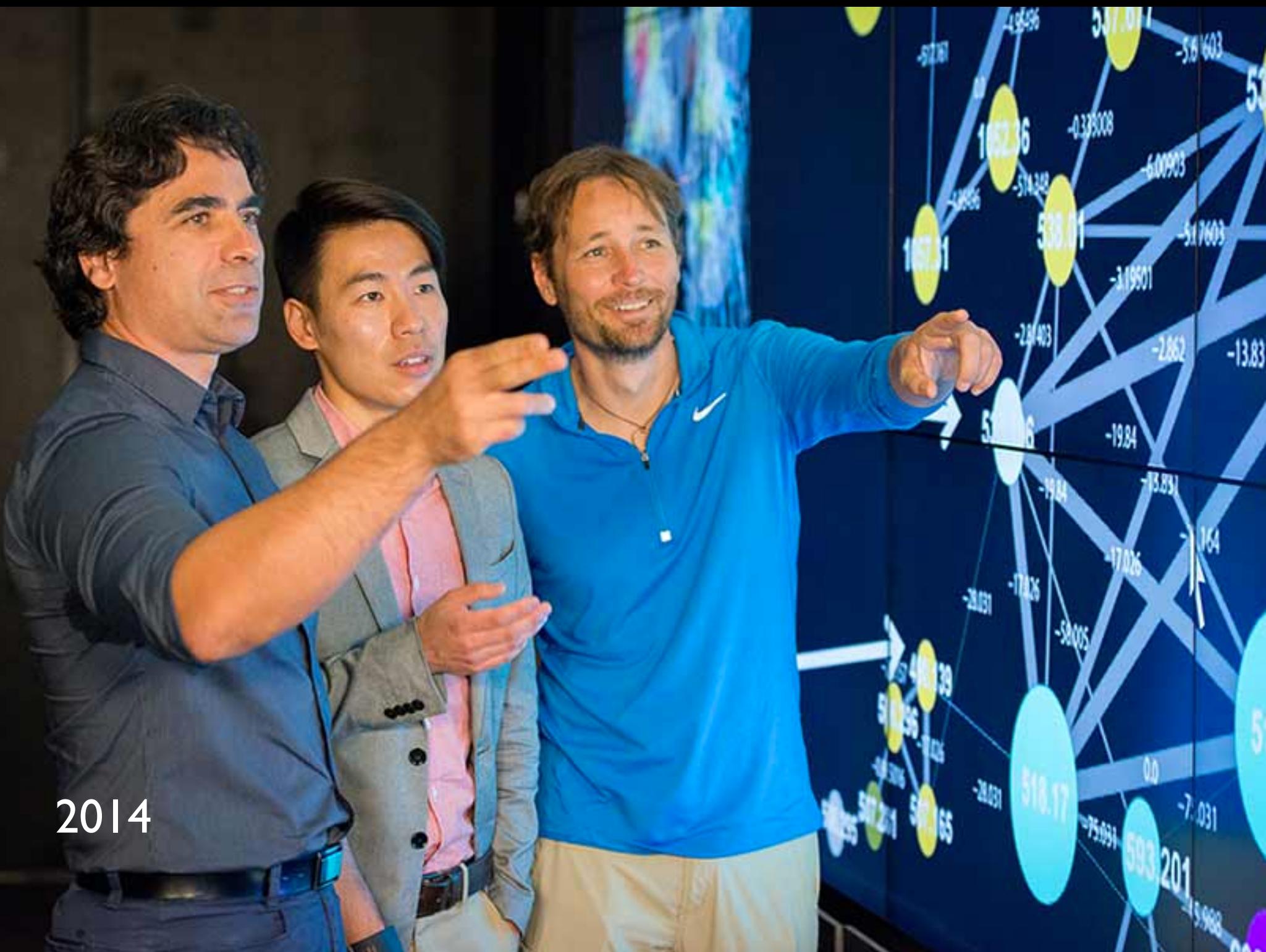
Where is this molecule  
found? any analogs?

## Status for small molecule mass spectrometry analysis

- 1) No centralized retrieval infrastructure
- 2) Not enough data in the public domain



2014



Being able to scale analysis is important.

>1,000 MS data files normal, must be possible to process  
100,000s to  
1,000,000s of files but ultimately billions of files.  
(currently almost 500,000 files in our infrastructure)

Scaling enables to disseminate knowledge to the world.



Written with 127 authors,

nature  
biotechnology

<http://www.nature.com/nbt/journal/v34/n8/full/nbt.3597.html>

## Sharing and community curation of mass spectrometry data with Global Natural Products Social Molecular Networking

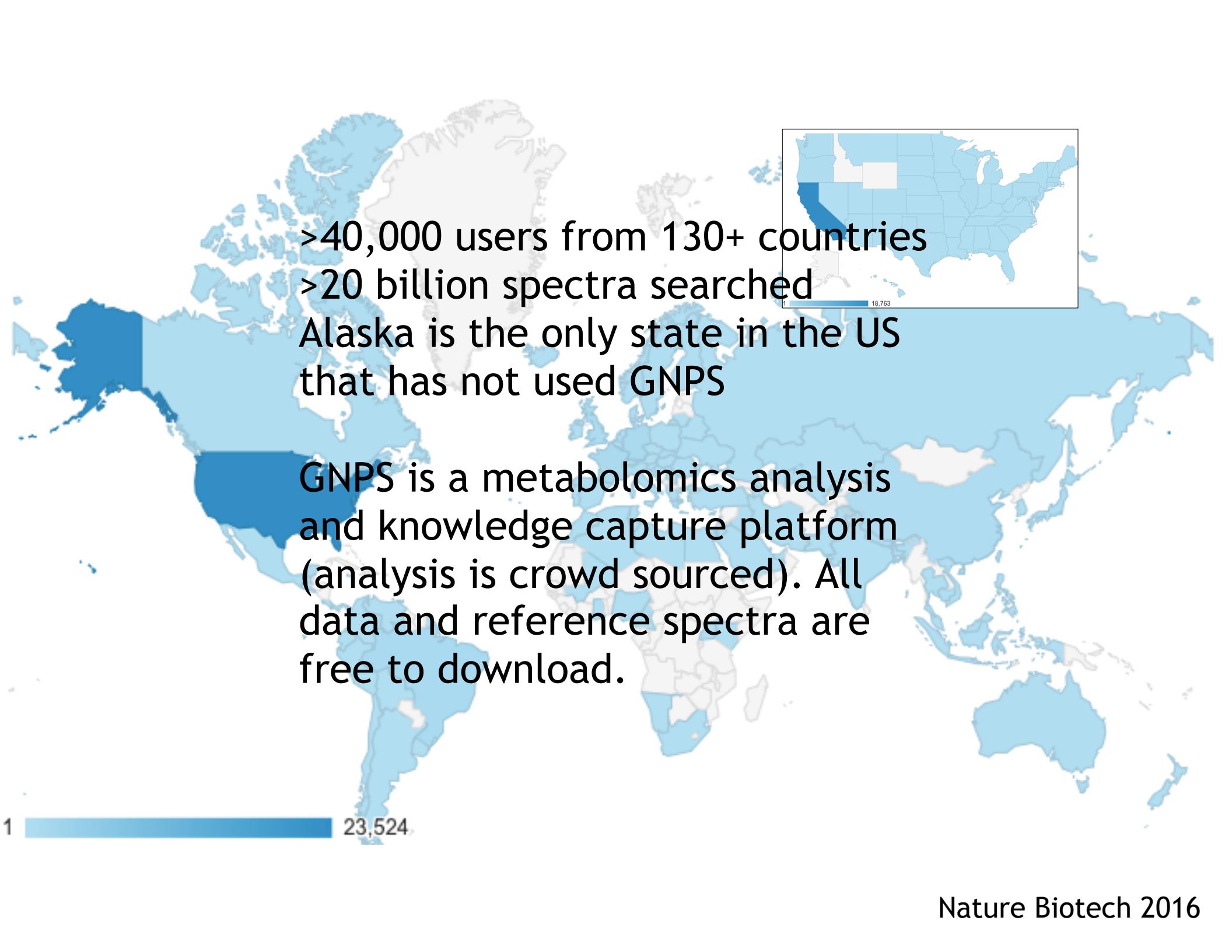
The potential of the diverse chemistries present in natural products (NP) for biotechnology and medicine remains untapped because NP databases are not searchable with raw data and the NP community has no way to share data other than in published papers. Although mass spectrometry (MS) techniques are well-suited to high-throughput characterization of NP, there is a pressing need for an infrastructure to enable sharing and curation of data. We present Global Natural Products Social Molecular Networking (GNPS; <http://gnps.ucsd.edu>), an open-access knowledge base for community-wide organization and sharing of raw, processed or identified tandem mass (MS/MS) spectrometry data. In GNPS, crowdsourced curation of freely available community-wide reference MS libraries will underpin improved annotations. Data-driven social-networking should facilitate identification of spectra and foster collaborations. We also introduce the concept of 'living data' through continuous reanalysis of deposited data.

NP from marine and terrestrial environments, including their inhabiting microorganisms, plants, animals, and humans, are routinely analyzed using MS. However, a single MS experiment can collect thousands of MS/MS spectra in minutes<sup>1</sup>, and individual projects can acquire millions of spectra. These data sets are too large for manual analysis. Furthermore, comprehensive software and proper computational infra-

and UniProt KnowledgeBase (UniProtKB), which provide robust platforms for data sharing and knowledge dissemination<sup>9,10</sup>. Recognizing the need for an analogous community platform to analyze NP MS data, we present GNPS. GNPS is a data-driven platform for the storage, analysis, and knowledge dissemination of MS/MS spectra that enables community sharing of raw spectra, continuous annotation of deposited data, and collaborative curation of reference spectra (referred to as spectral libraries) and experimental data (organized as data sets).

GNPS provides the ability to analyze a data set and to compare it to all publicly available data. By building on the computational infrastructure of the University of California San Diego (UCSD) Center for Computational Mass Spectrometry (CCMS; <http://proteomics.ucsd.edu/>), GNPS provides public data set deposition and/or retrieval through the Mass Spectrometry Interactive Virtual Environment (MS-IVE) 1.0 interface. The GNPS platform is currently

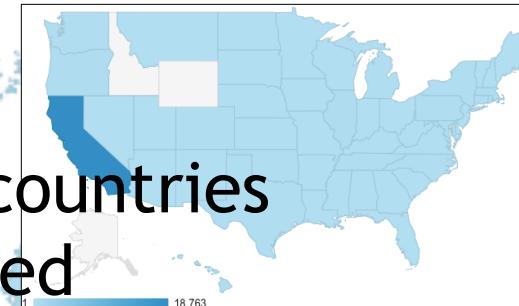


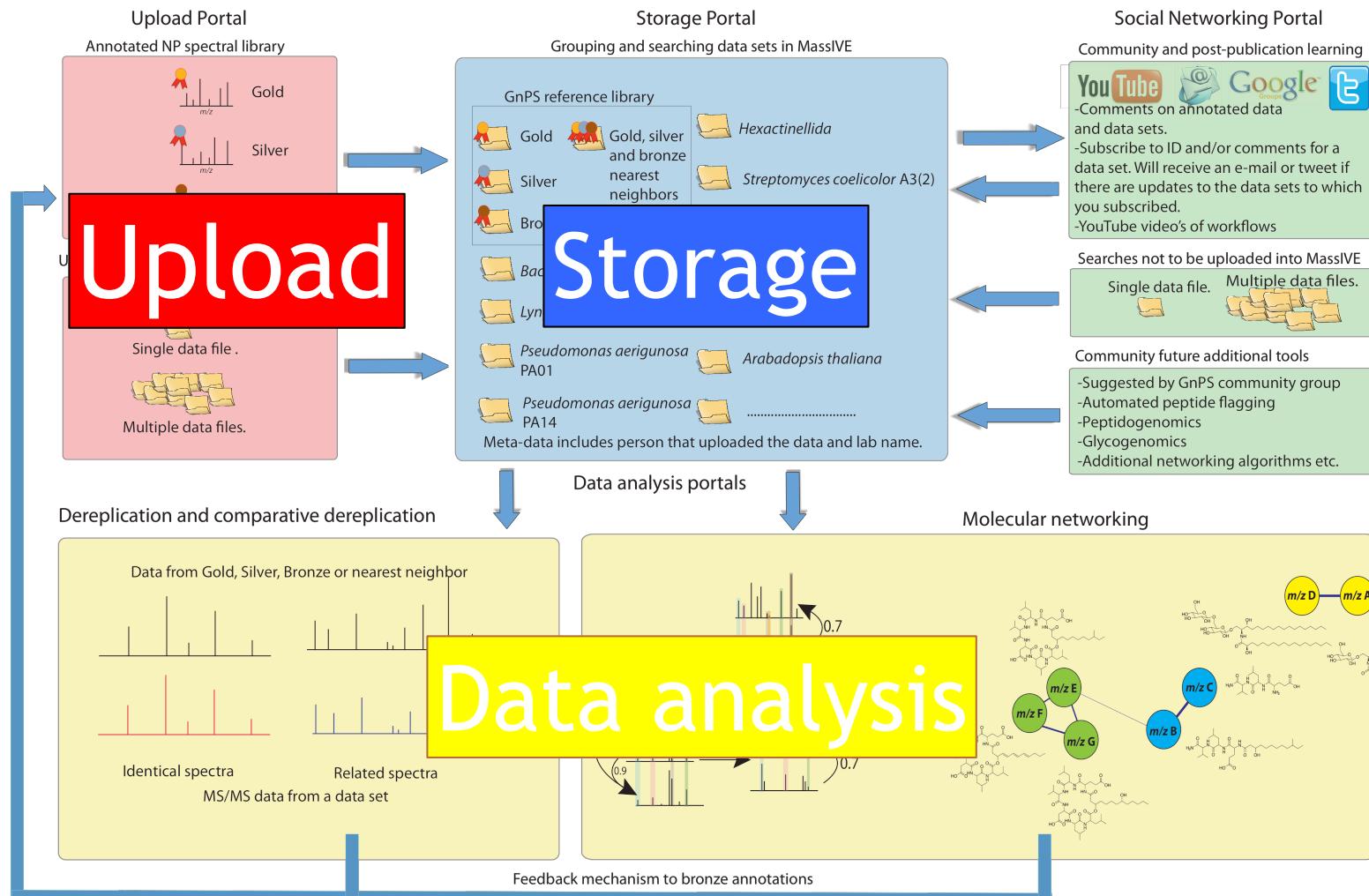


>40,000 users from 130+ countries  
>20 billion spectra searched  
Alaska is the only state in the US  
that has not used GNPS

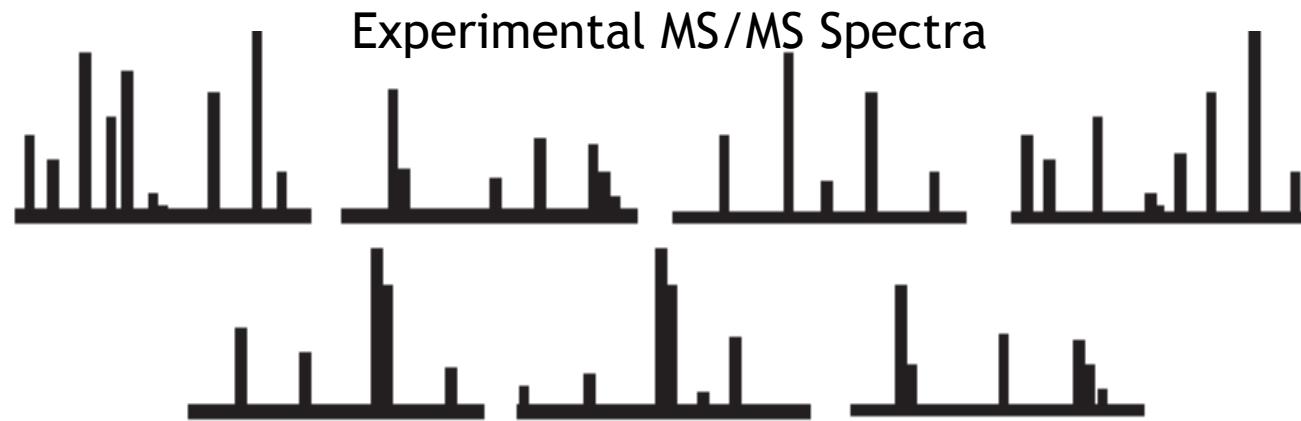
GNPS is a metabolomics analysis  
and knowledge capture platform  
(analysis is crowd sourced). All  
data and reference spectra are  
free to download.

23,524

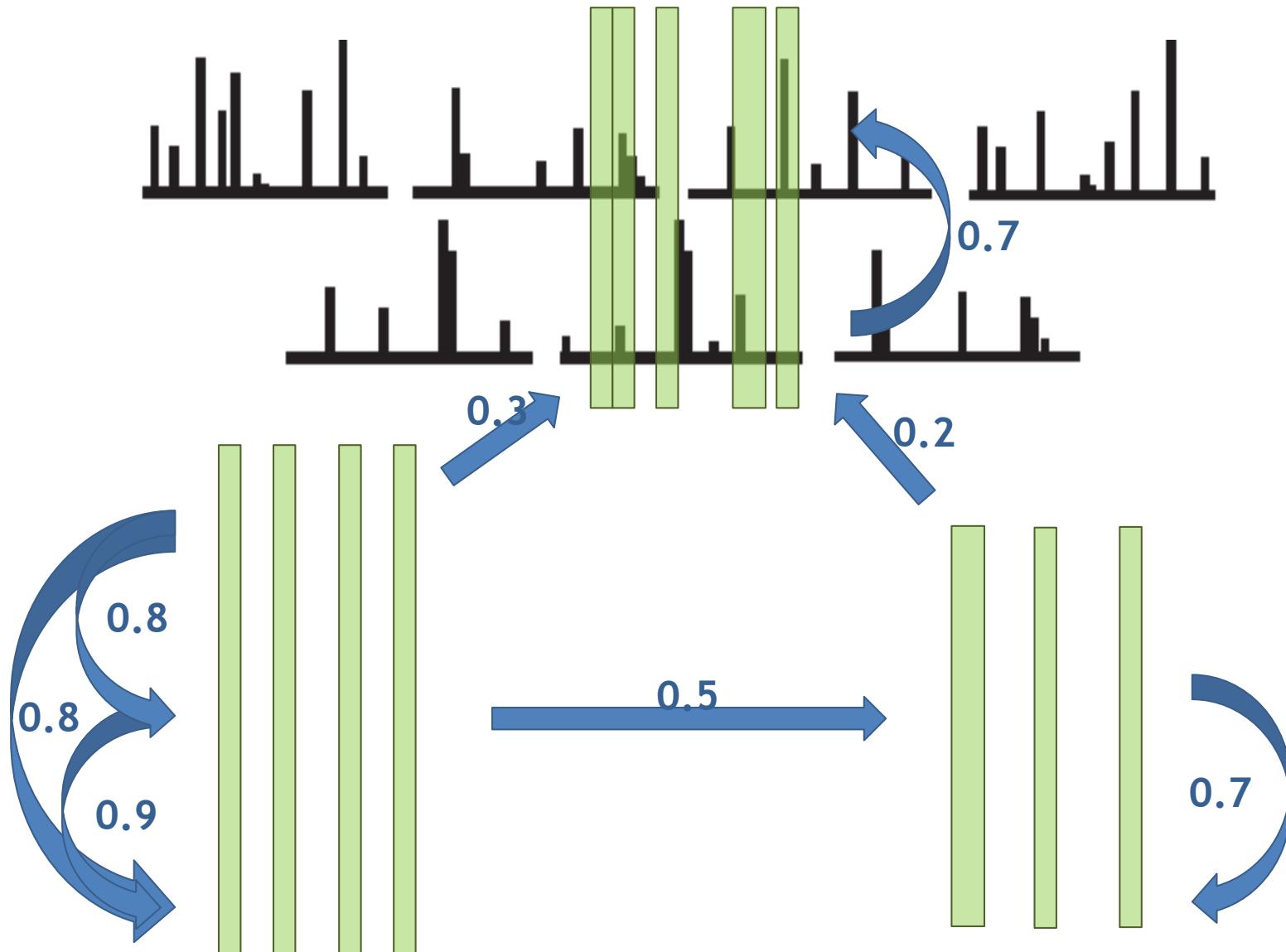




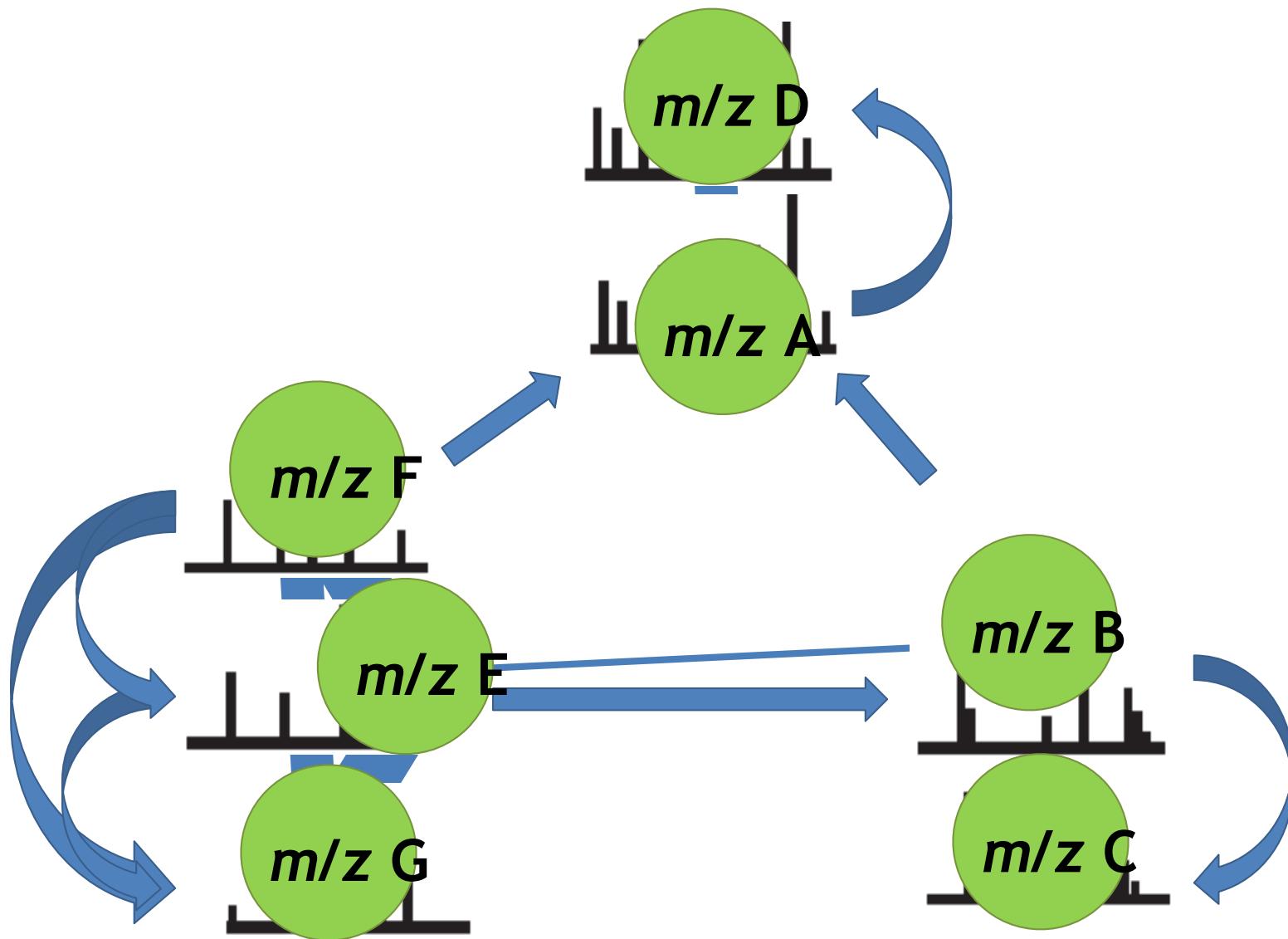
## Molecular networking



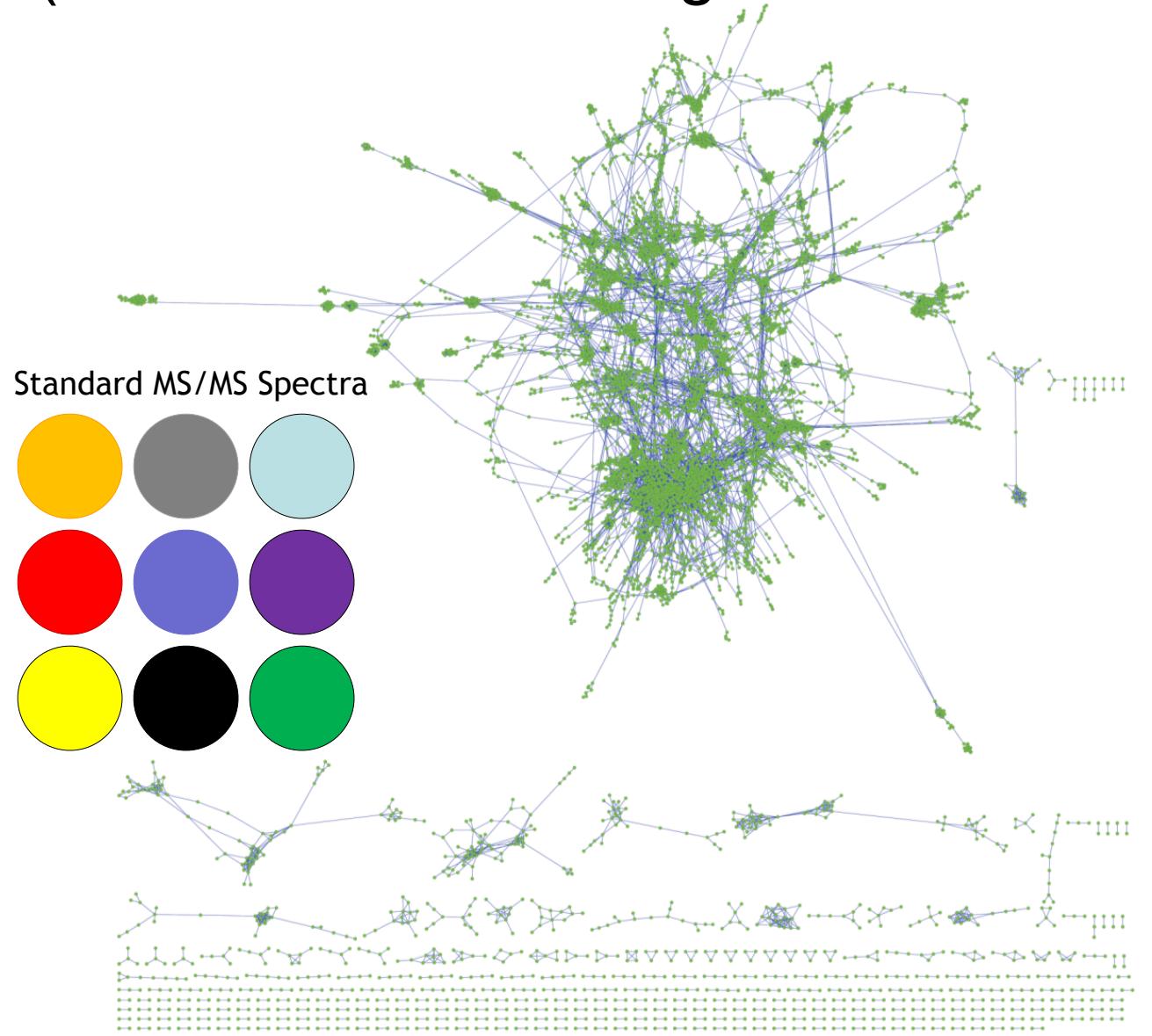
## Molecular networking

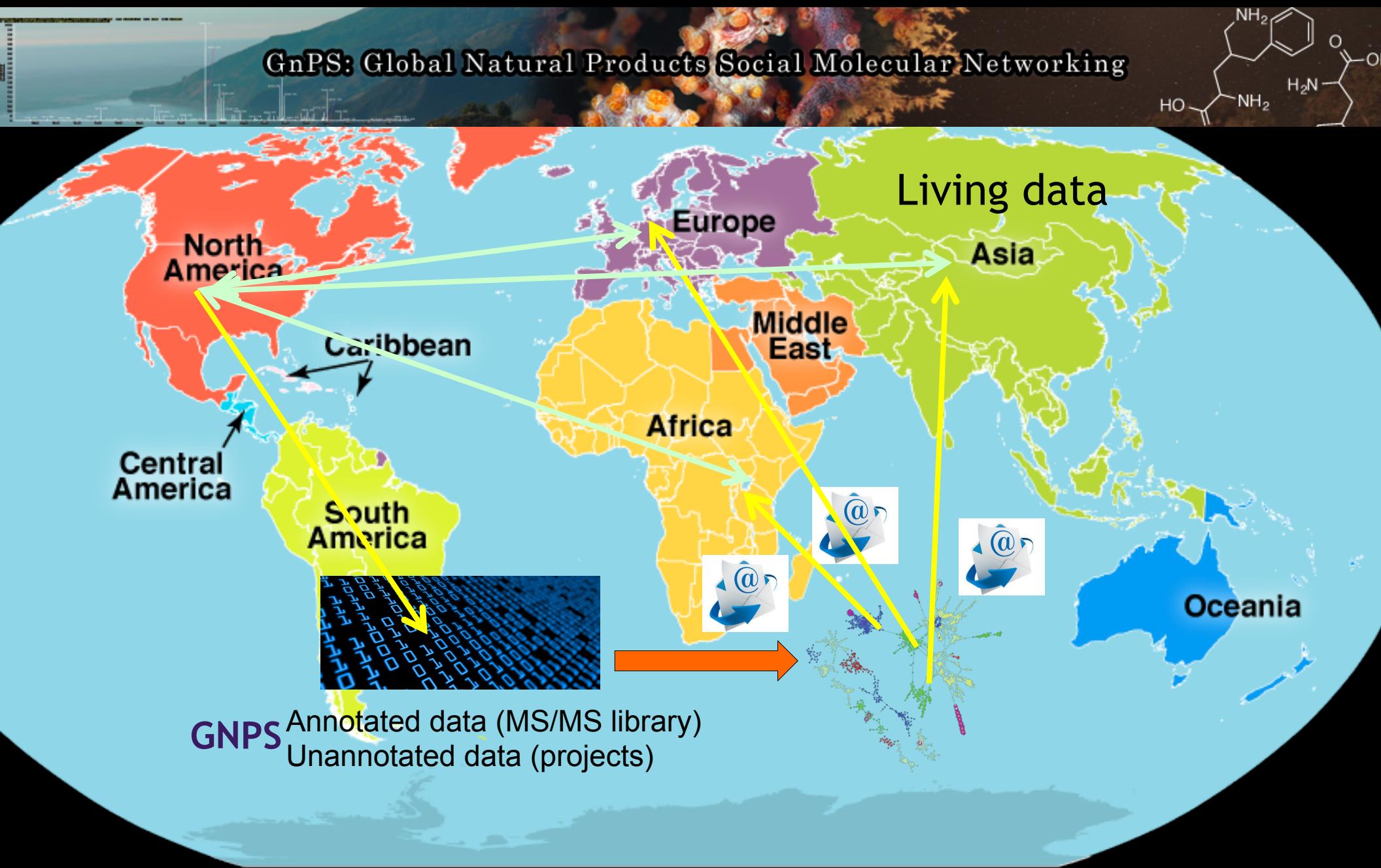


## Molecular networking

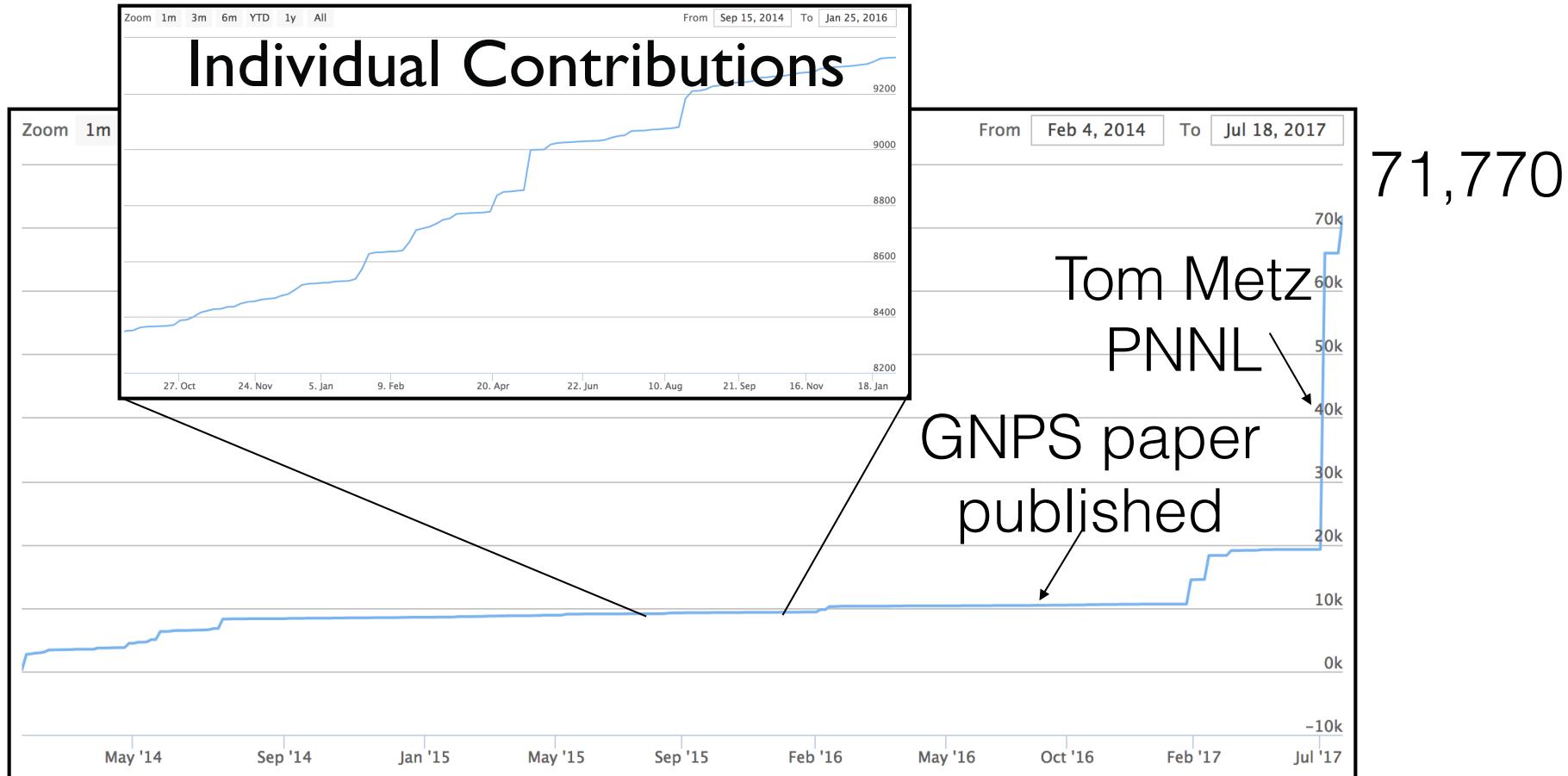


# Dereplication and comparative dereplication (also known as finding known unknowns)





# Growth of knowledge in GNPS (contributed by the GNPS community of >40,000 people)



# Data analysis reproducibility



GNPS: Global Natural Products Social Molecular Networking

Logout | My User | Update Profile | Jobs | MassIVE Datasets | Documentation | Forum | Contact

NCC(CN)C(=O)C[C@H](N)C[C@H](O)CN1CCC[C@H]1c2ccccc2N

[Back to main page](#)

ProteoSAFe Workflow Tasks

◀ Hits 1 ~ 30 out of 329

▶ Go to

Go

Select columns

	Description	User	Workflow	Site	Status	Create Time	Elapsed Time
1	cheese ID=a9bef02315e546949f3c59181b4eaa33	pderrestein	METABOLOMICS-SNETS	GNPS	DONE	Nov. 2, 2017, 4:48 PM	9 day 23:38:00
2	fungi vs food ID=8254edfe3528434a93dcda0419f7f275	pderrestein	METABOLOMICS-SNETS	GNPS	DONE	Jul. 28, 2017, 8:04 AM	1 day 04:17:55
3	food vs skin ID=19d3b6b5f3864fb3b0976ee400379e12	pderrestein	METABOLOMICS-SNETS	GNPS	DONE	Jul. 28, 2017, 8:02 AM	0 day 22:58:11
4	lululemon 5 ion cos 0.7 min cluster size 3 parent mass 0.1 ans MS/mS match 0.1 ID=8f1b79bcc4a4478972f22bf4f8491bf	pderrestein	METABOLOMICS-SNETS	GNPS	DONE	Jul. 26, 2017, 8:54 PM	1 day 18:55:49
5	lululemon dereplication ID=69422dc3ac6c45318ff395caeff21b8	pderrestein	MOLECULAR-LIBRARYSEARCH	GNPS	DONE	Jul. 26, 2017, 8:47 AM	0 day 03:55:49
6	lululemon dereplication ID=4b2b0712650f4807af6e37450b905250	pderrestein	MOLECULAR-LIBRARYSEARCH	GNPS	FAILED	Jul. 26, 2017, 8:46 AM	0 day 00:03:53
7	ID=958b362704ee4bfeb63d3c0a71c8906e	pderrestein	UPDATE-SINGLE-ANNOTATED-BRONZE	GNPS	DONE	Jul. 24, 2017, 5:59 PM	0 day 00:01:21

# Data analysis reproducibility

GNPS: Global Natural Products Social Molecular Networking

[Logout](#) | [My User](#) | [Update Profile](#) | [Jobs](#) | [MassIVE Datasets](#) | [Documentation](#) | [Forum](#) | [Contact](#)

NCC(CN[C@@H](C)CO)C1CCC(N)C(Cc2ccccc2)CC1

## METABOLOMICS-SNETS

DONE

[Clone]

[Restart][Delete]

[ [View All Library Hits](#) | [View All Clusters With IDs](#) | [View All Compounds](#) ]

### 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](#) | [View Compounds and File Occurrence](#) ]

### Advanced Views - Networking Graphs

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

### Advanced Views - Third Party Visualization

[ [View Emporer PCoA Plot in GNPS](#) | [View ili in GNPS](#) ]

### Community Matches

[ [Dataset Matches](#) ]

### Network Visualizations

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

### Export

[ [Download Clustered Spectra as MGF](#) | [Download Cytoscape Data](#) | [Download Bucket Table](#) | [Download BioM For Qiime/Qiita](#) | [Download Metadata For Qiime](#) | [Make Public Dataset](#) | [Download ili Data](#) ]

# Data analysis reproducibility and role of digitizing metadata

The screenshot shows a web browser window with the following details:

- GNPS: Global Natural Products Social Molecular Networking** is displayed at the top.
- The URL in the address bar is <https://hsmail.ucsd.edu/owa/>.
- The browser toolbar includes various icons for file operations (New, Open, Save, Print, PDF) and links to other services like Google, Twitter, and LinkedIn.
- The main content area shows an email inbox titled "Inbox" with 1000 results. A search bar at the top of the inbox lists "julia".
- The inbox list shows several messages from different senders:

  - Julia Gauglitz: Re: could you send me the slides you had today, I want to p... (6:13 PM)
  - Dorrestein, Pieter: could you send me the slides you had today, I want to partic... (4:58 PM)
  - Sandrine Miller-Montgomery: Re: 2018 Seed Grant Plans (4:00 PM)
  - Dorrestein, Pieter: FW: cheese network (3:08 PM)
  - Julia Gauglitz: cheese network (2:32 PM)
  - Knight, Rob: Re: 2018 Seed Grant Plans (2:21 PM)
  - Knight, Rob: Re: 2018 Seed Grant Plans (2:17 PM)
  - Austin Swafford: Re: 2018 Seed Grant Plans (2:10 PM)
  - Daniel Freed: Re: 2018 Seed Grant Plans (1:58 PM)
  - Kelly C Weldon: Monthly Update-April 2018 (1:57 PM)

- To the right of the inbox, a message is open from **Julia Gauglitz** [julia.gauglitz@gmail.com] with the subject **cheese network**. The message body contains the URL <https://gnps.ucsd.edu/ProteoSAFe/status.jsp?task=4a637ad8d8794d4a9d416e8ba637ad32>.
- The message header shows the recipient as **Dorrestein, Pieter** and the cc as **Elizabeth Brown [eabrown@ucsd.edu]**.
- The timestamp for the message is **Tuesday, May 01, 2018 2:31 PM**.
- The message notes that it was forwarded from the original sender on **5/1/2018 3:08 PM**.
- The bottom left of the screen shows a sidebar with icons for Mail, Calendar, Contacts, Tasks, and Public Folders.

## Metadata associated with this data set

Not yet done but IMO we need to capture metadata at the time of deciding to run the experiment, not add afterwards.

[https://gnps.ucsd.edu/ProteoSAFe/status.jsp?  
task=4a637ad8d8794d4a9d416e8ba637ad32](https://gnps.ucsd.edu/ProteoSAFe/status.jsp?task=4a637ad8d8794d4a9d416e8ba637ad32)

# nature protocols

Recipes for Researchers

[www.natureprotocols.com](http://www.natureprotocols.com)

GNPS: Global

## METABOLOMICS-SNETS

DONE

[Clone]

[ View All Library Hits | View ]

Methods and Citation for Ma

[ Networking Parameters and ]

Experimental Views

[ Reanalyze Cluster Spectra ]

Auxiliary Views

[ View Network, Node Centric ]

Advanced Views - Networkin

[ Nodes, MZ Histogram | Edge ]

Advanced Views - Third Part

[ View Emporer PCoA Plot in ]

Community Matches

[ Dataset Matches ]

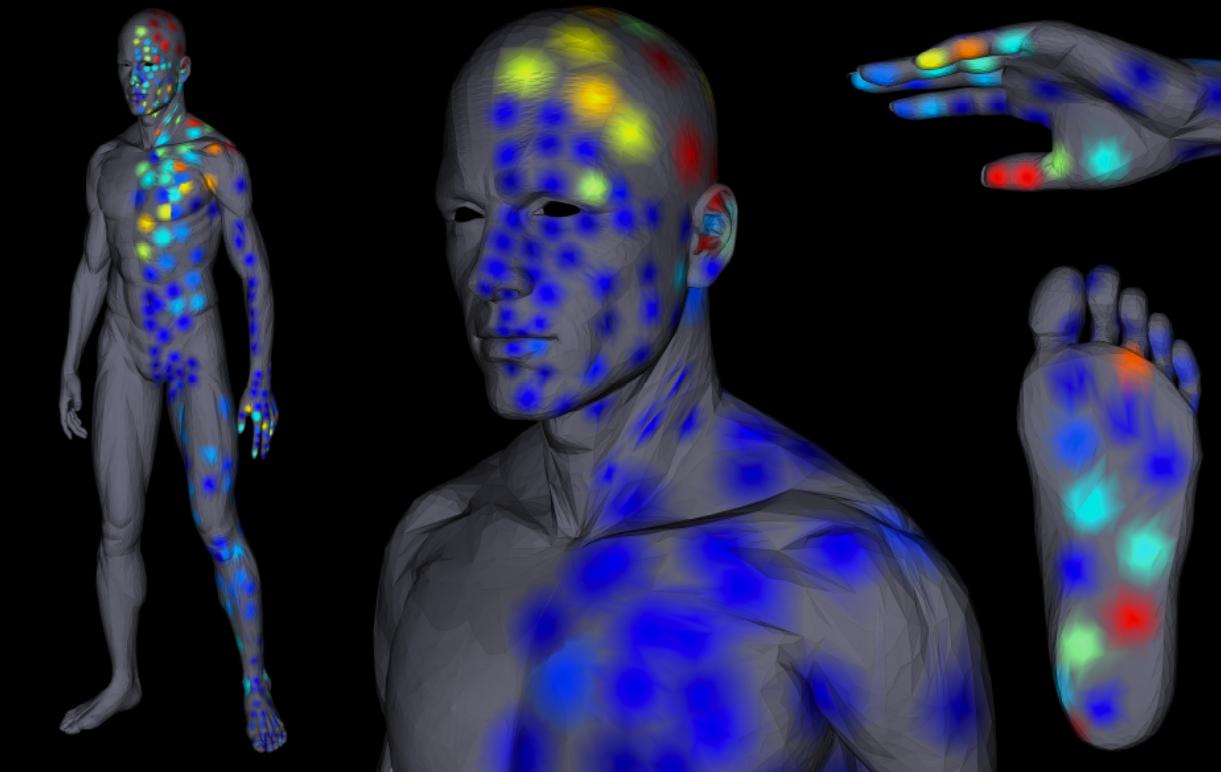
Network Visualizations

[ View Spectral Families (In B ]

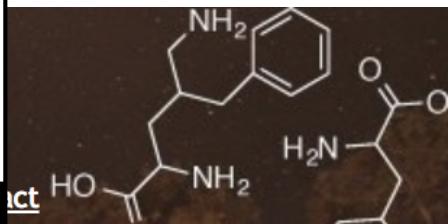
Export

[ Download Clustered Spectra ]

Metadata For Qiime | Make Pe



- » **Mouse cerebellar window**
- » **HPLC detection of monoamines and their cofactors**
- » **Probing virus–host interactions with AFM**
- » **Massively parallel cancer-mutation stratification**
- » **Chemical synthesis of histones**



[ Restart ] [ Delete ]

[ Reference ]

[ MZ Histogram ]

[ For Qiime/Qiita | Download ]

Examples of data reuse  
from my own lab.

# Example Data reuse 1.

1    **Bioactivity Based Molecular Networking for the Discovery of Active Natural**  
2    **Products**

3  
4    Louis-Félix Nothias,<sup>1,2,3</sup> Mélissa Nothias-Esposito,<sup>3,4</sup> Ricardo da Silva,<sup>1,2</sup> Ivan Protsyuk,<sup>5</sup> Pieter  
5    Leyssen,<sup>6</sup> David Touboul,<sup>3</sup> Julien Paolini,<sup>4</sup> Theodore Alexandrov,<sup>5</sup> Marc Litaudon,<sup>3</sup> & Pieter C.  
6    Dorrestein<sup>1,2,\*\*</sup>

7  
8  
9  
10    **ABSTRACT**

11    It is a common problem in natural product therapeutic lead discovery programs that  
12    despite good bioassay results in the initial extract, the active compound(s) could not be  
13    isolated during subsequent bioassay-guided purification. Herein, we present the  
14    concept of bioactive molecular networking to find the candidate active molecules directly  
15    from fractionated bioactive extracts. Mass spectrometry based molecular networking,  
16    bioactivity score prediction to create bioactive molecular networks enables the discovery  
17    and the targeted isolation of lead molecules. We highlight the feasibility of this approach

# Why study Euphorbia dendroides for anti-viral activities?



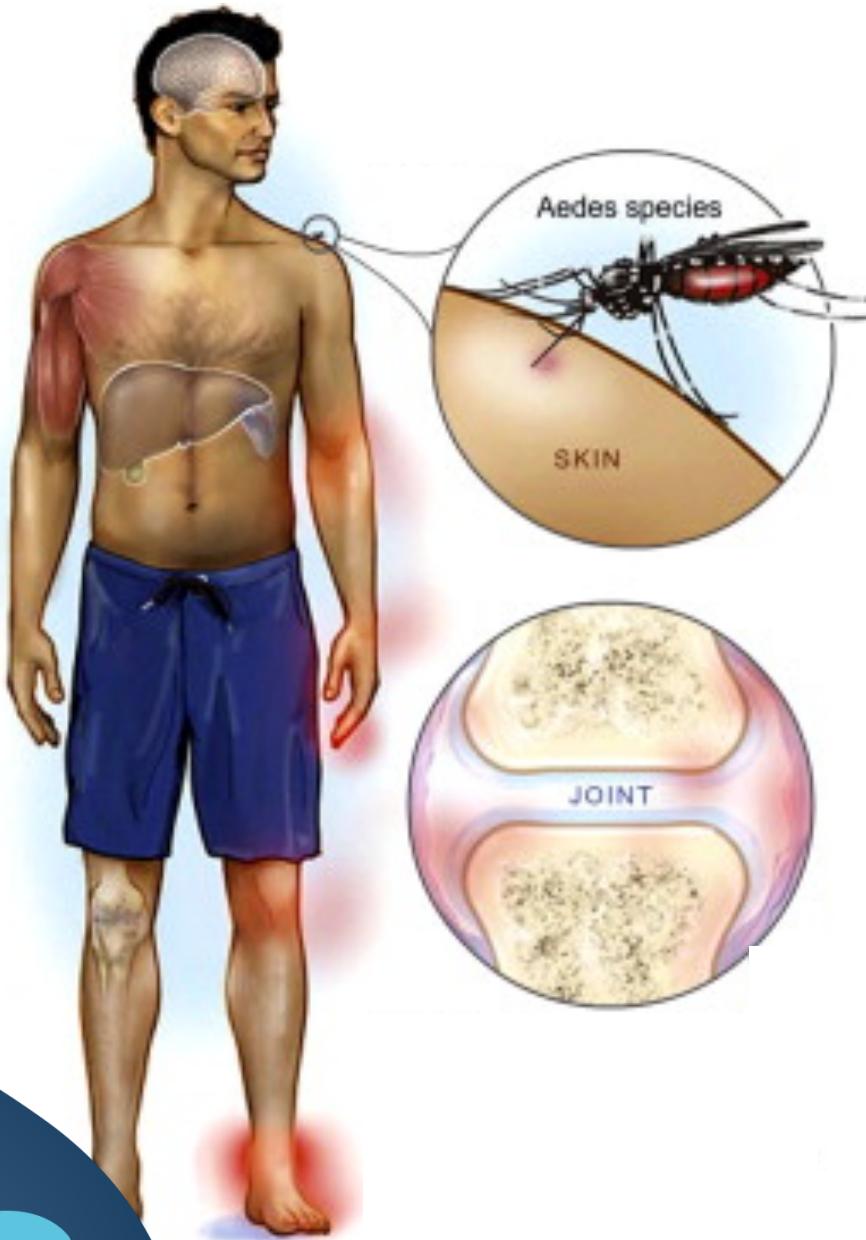
# Why study Euphorbia dendroides for anti-viral activities?



Warts=Viral

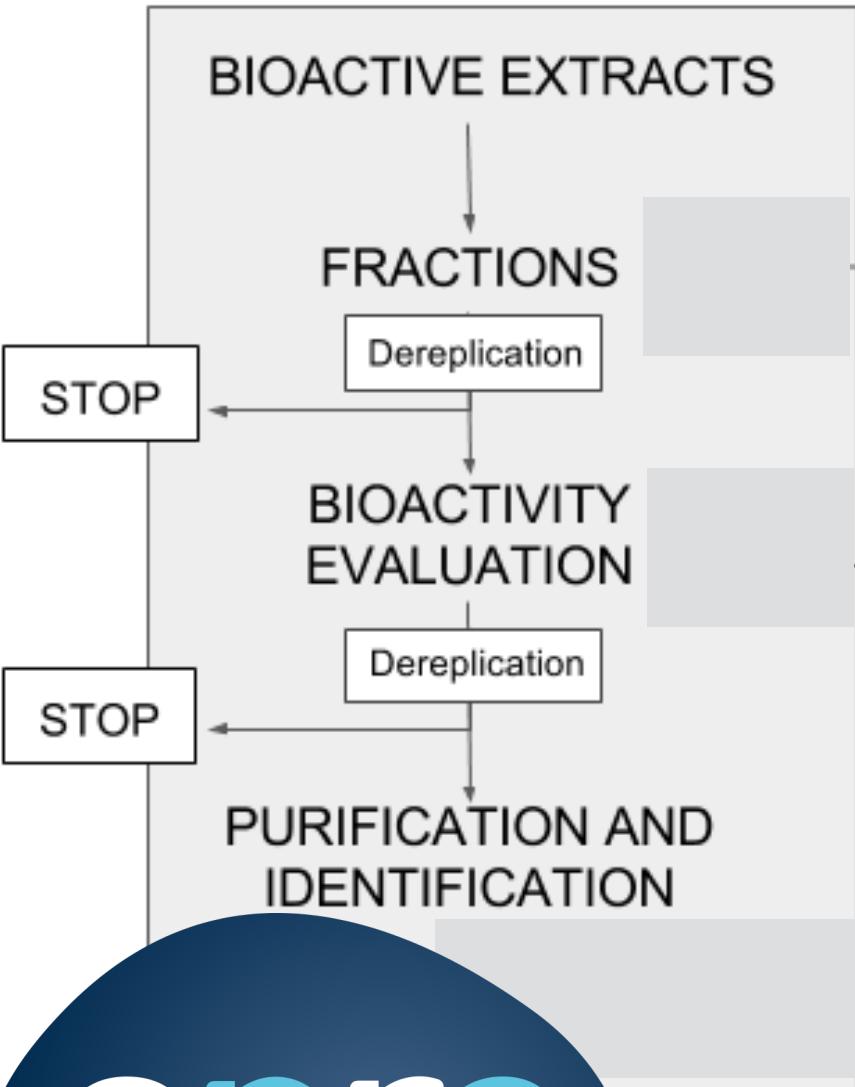
Marc Litaudon and Julien Paolini asked - Does the extract of this plant have anti-viral activity?





Tested extracts against chikungunya virus  
in monkey cell lines

## “Classical” bioassay-guided fractionation



Euphorbia dendroides

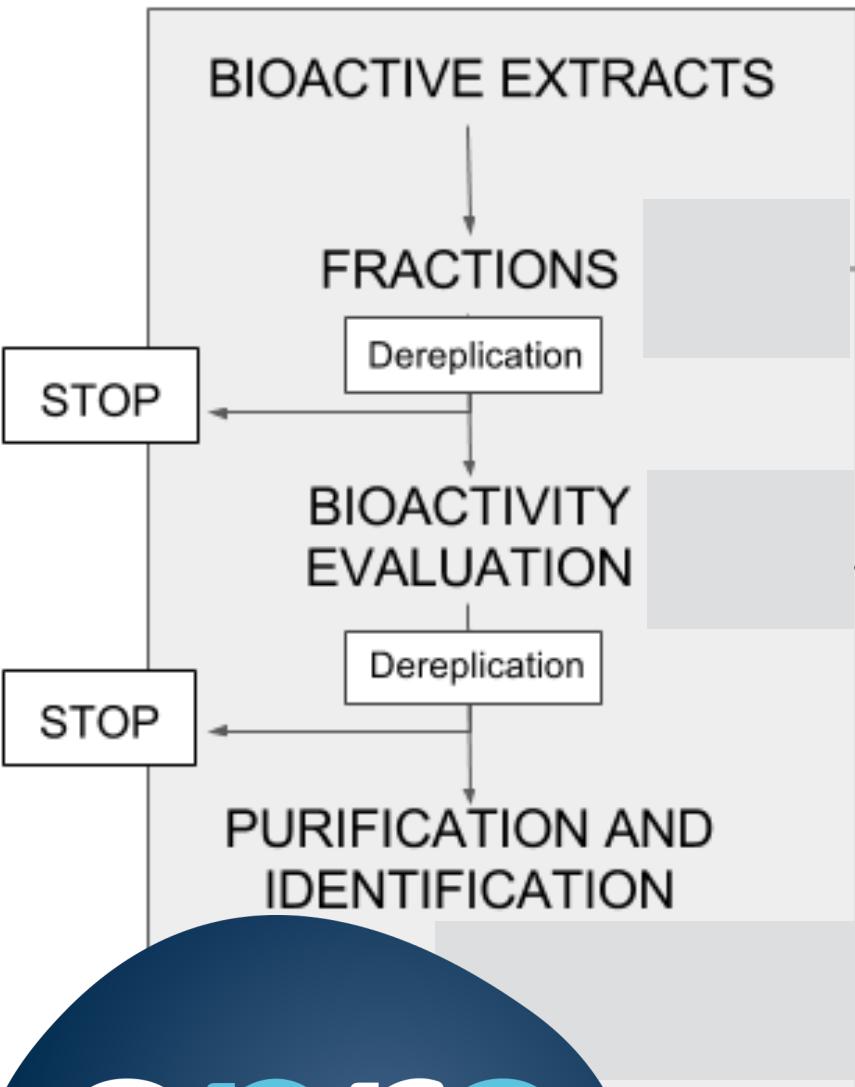
Activity against the chikungunya virus

Fractionated into 20 fractions

3 fractions active EC<sub>50</sub>= 0.268,  
0.169 and 1.130 µg/ml



## “Classical” bioassay-guided fractionation

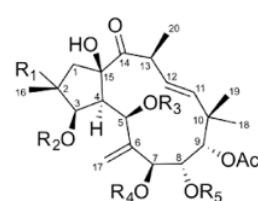


Activity against the chikungunya virus

Fractionated into 20 fractions

3 fractions active EC<sub>50</sub>= 0.268,  
0.169 and 1.130 µg/ml

**PREVIOUSLY ISOLATED COMPOUNDS (1-19)**



III  
I  
III

R<sub>1</sub>

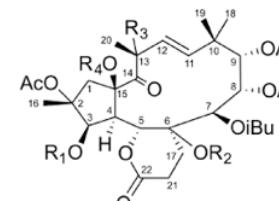
R<sub>2</sub>

R<sub>3</sub>

R<sub>4</sub>

R<sub>5</sub>

II  
I  
III



R<sub>1</sub>

R<sub>2</sub>

R<sub>3</sub>

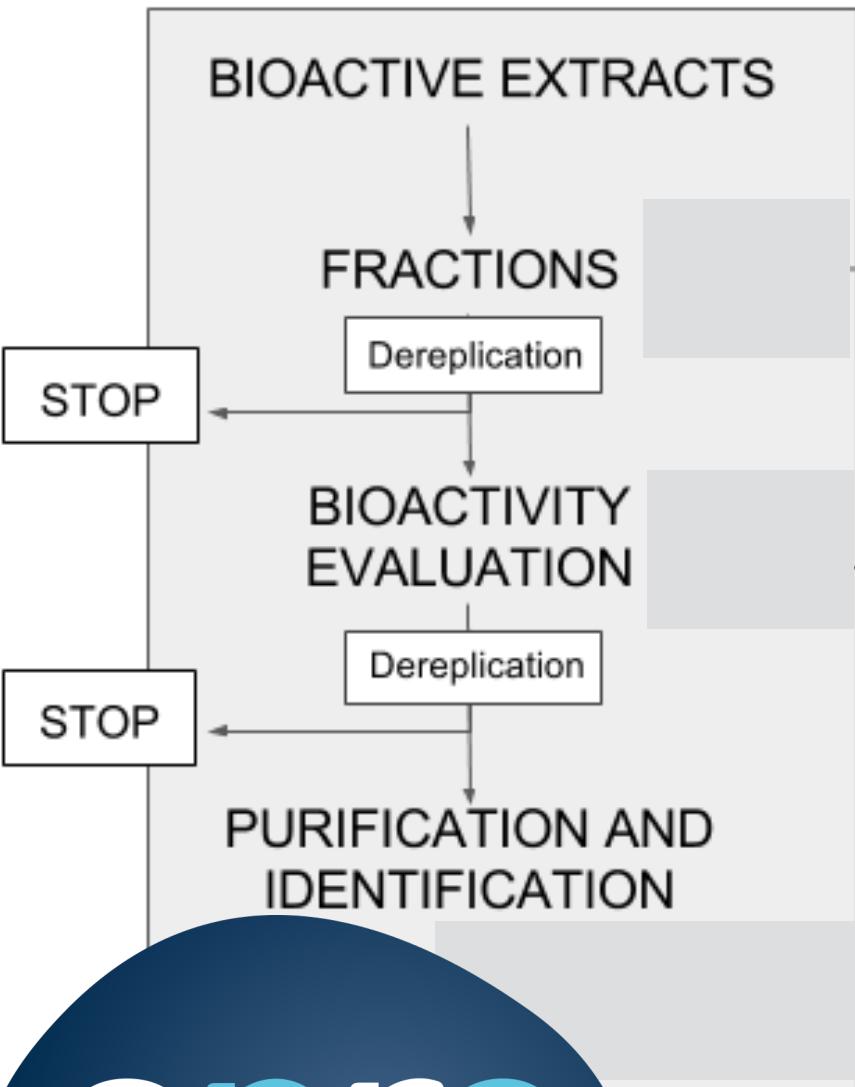
R<sub>4</sub>

R<sub>5</sub>

R<sub>1</sub>

</div

## “Classical” bioassay-guided fractionation

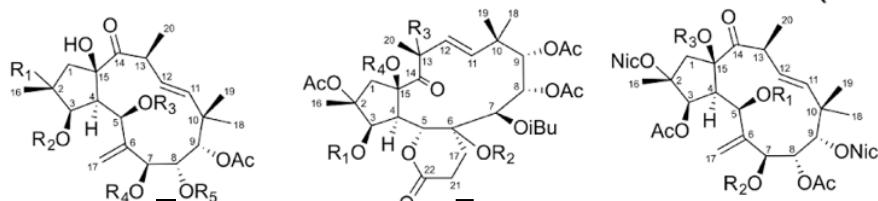


Activity against the chikungunya virus

Fractionated into 20 fractions

3 fractions active EC<sub>50</sub> = 0.268,  
0.169 and 1.130 µg/ml

**PREVIOUSLY ISOLATED COMPOUNDS (1-19)**



	R <sub>1</sub>	R <sub>2</sub>	R <sub>3</sub>	R <sub>4</sub>	R <sub>5</sub>														
1	H	Ac	iBu	Bz	Ac														
2	OH	Ac	iBu	Bz	Ac														
3	OH	Bz	H	Bz	Ac														
4	OAc	H	iBu	Bz	Ac														
5	iBu	Ac	iBu	Bz	Ac														
6	OAc	Ac	iBu	Bz	Ac														
7	OAc	A	Ac	iBu	Ac														
8	OAc	Bz	H	iBu	Ac														
9	OAc	H	iBu	Bz	Ac														
10	OAc	H	Bz	Bz	Ac														
11	OAc	H	Bz	Bz	H														
12																			
13																			
14																			
15																			
16																			
17																			
18																			
19																			

**Inactive-person years**

Compounds 1-19 did not show significant bioactivity against CHIKV replication [20]

A =

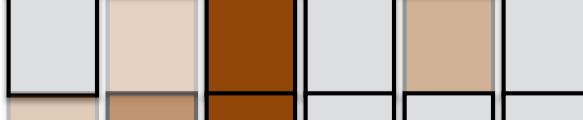
**Activity**  
(EC50 in  $\mu\text{g}/\text{ml}$ )

*m/z* A B C D E F

**Fraction 1**



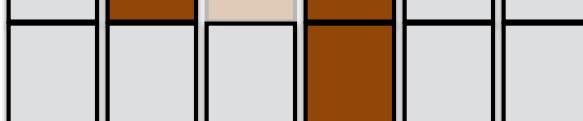
**Fraction 2**



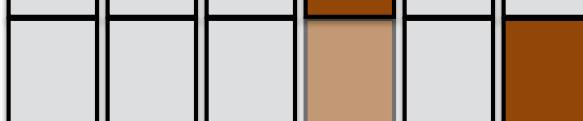
**Fraction 3**



**Fraction 4**



**Fraction 5**

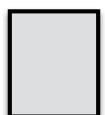
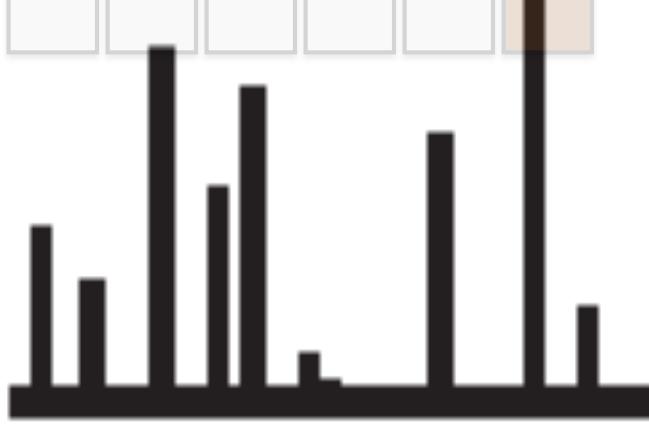


**Fraction 6**

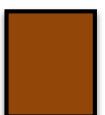


**Fract....**

**Etc.**



= not detected/no activity

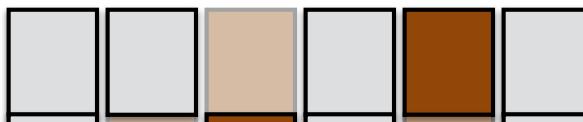


= most intense/active

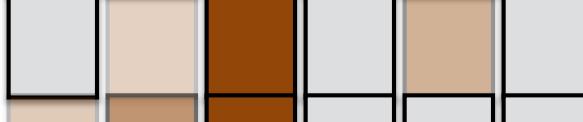
**Activity**  
(EC50 in  $\mu\text{g}/\text{ml}$ )

*m/z* A B C D E F

**Fraction 1**



**Fraction 2**



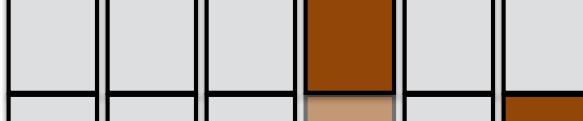
**Fraction 3**



**Fraction 4**



**Fraction 5**

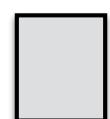
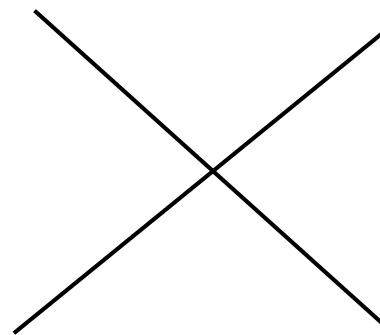
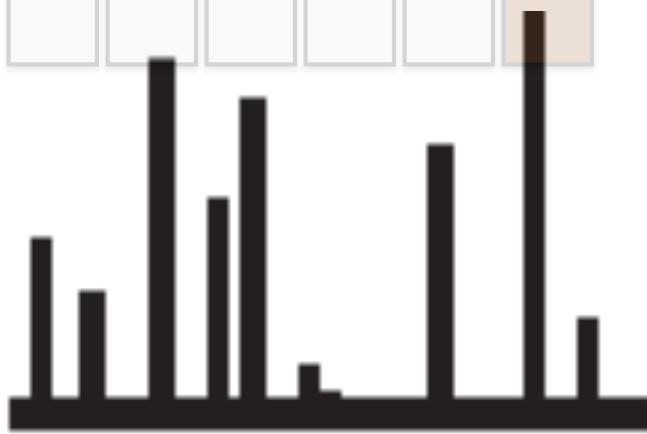


**Fraction 6**

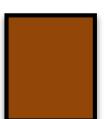


**Fract....**

**Etc.**

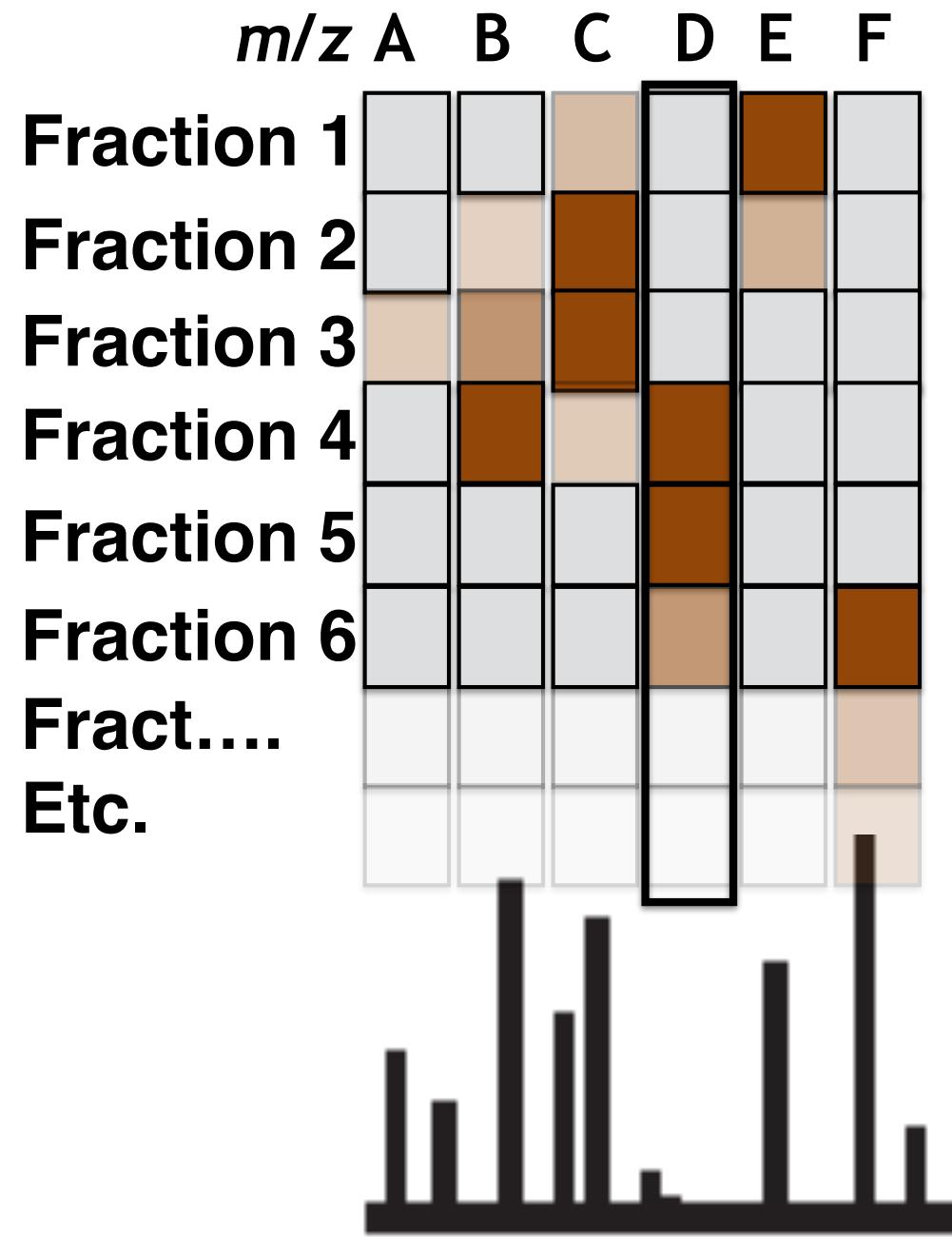


= not detected/no activity

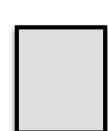
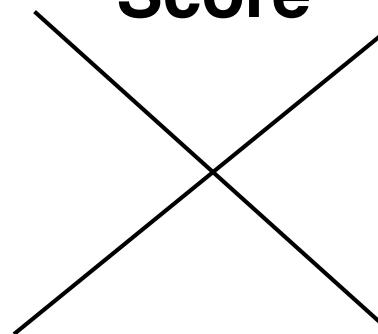


= most intense/active

**Activity**  
(EC50 in  $\mu\text{g}/\text{ml}$ )

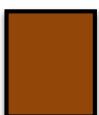


**Pearson  
Score**



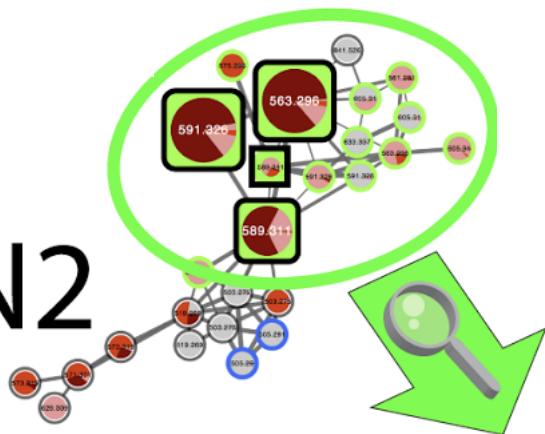
*m/z*

= not detected/no activity



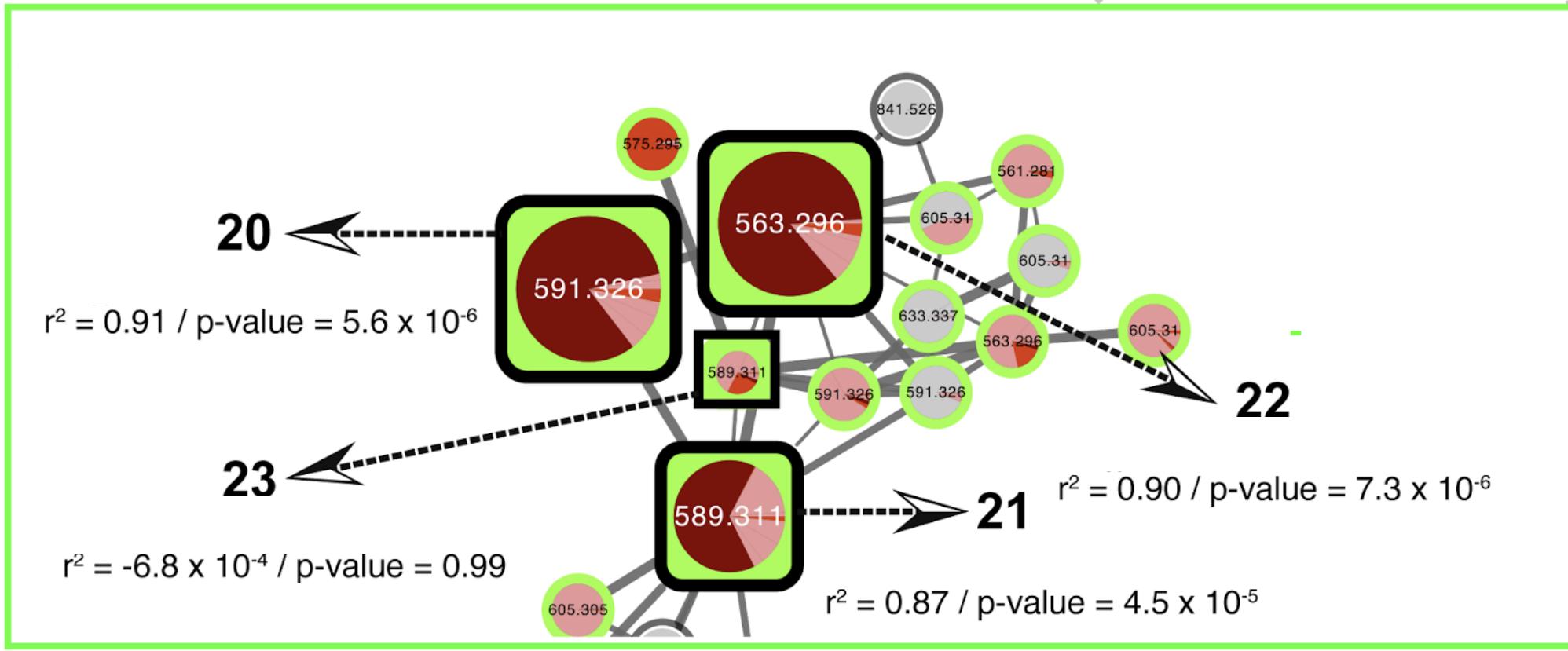
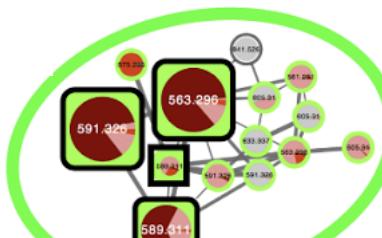
= most intense/active

# MN2



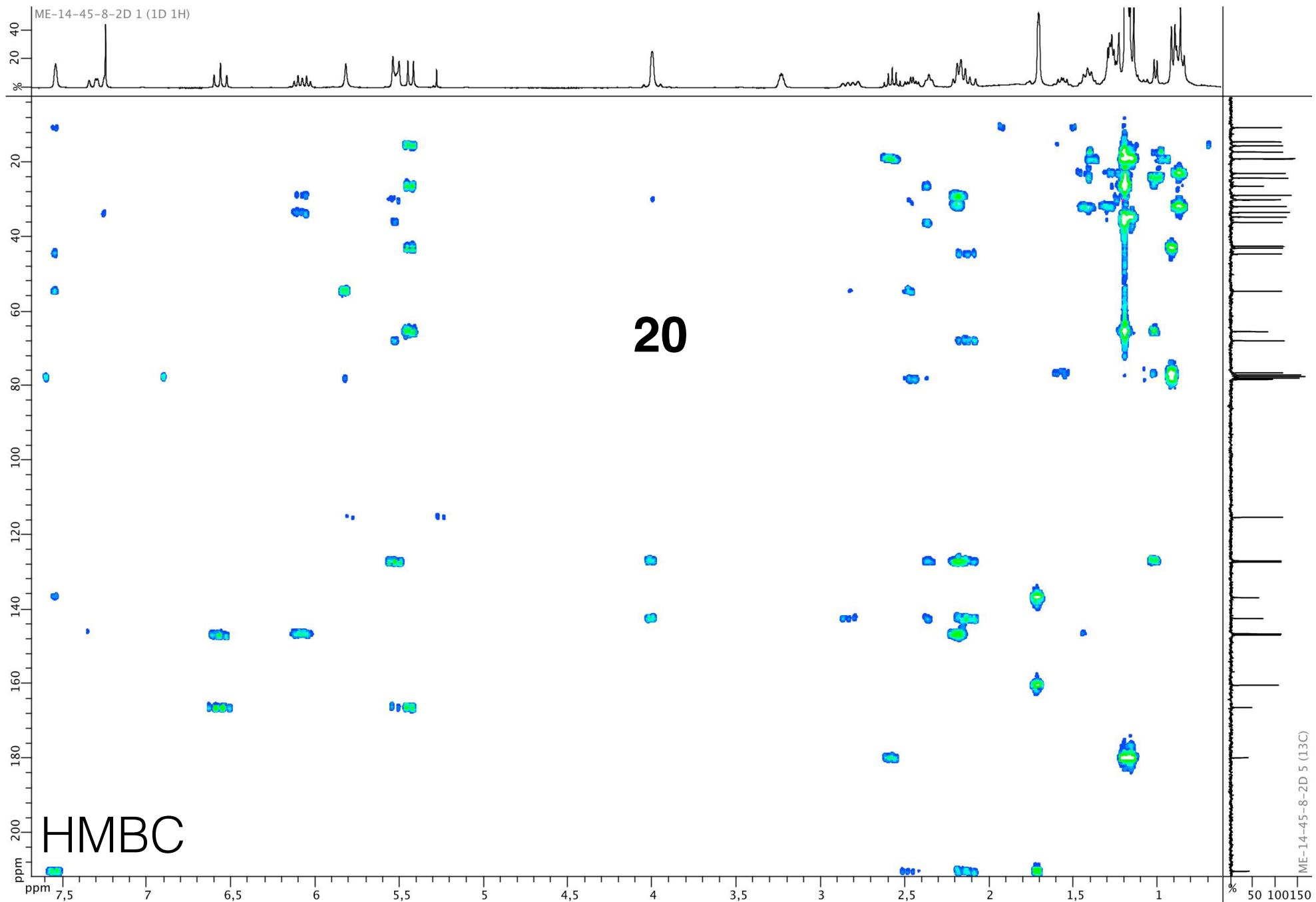
Correlation does not mean causation!!  
We have potential NP candidates responsible for activity.

# MN2

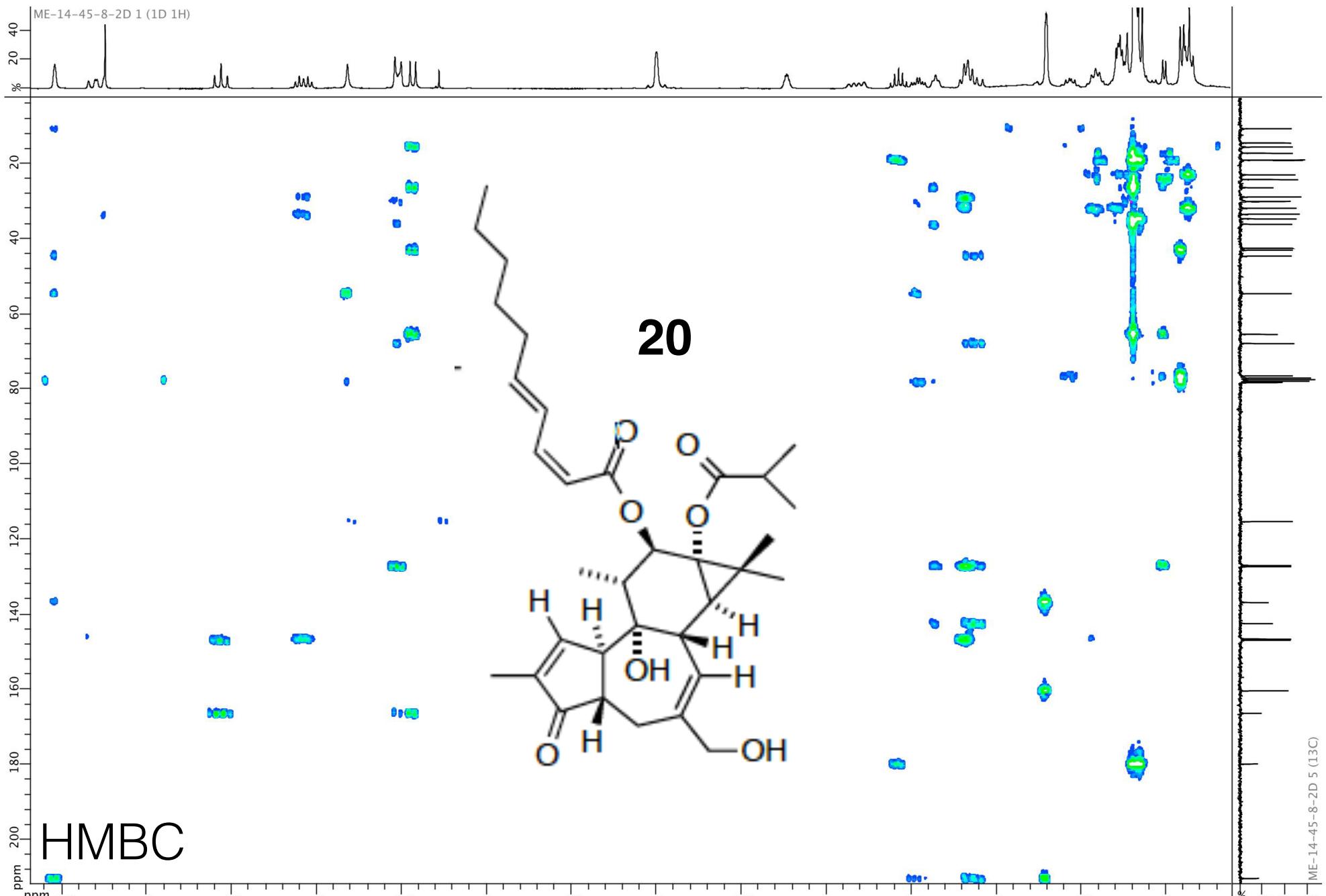


Correlation does not mean causation!!  
We have potential NP candidates responsible for activity.

# MS based targeted isolation - NMR



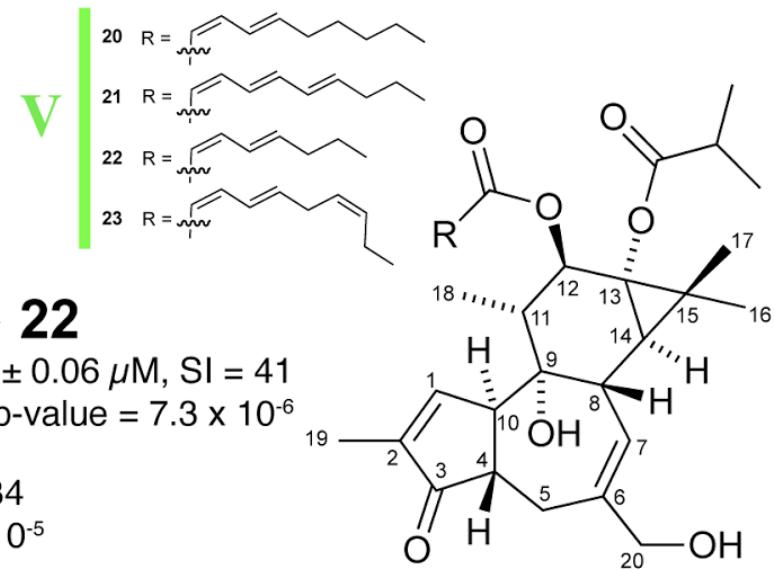
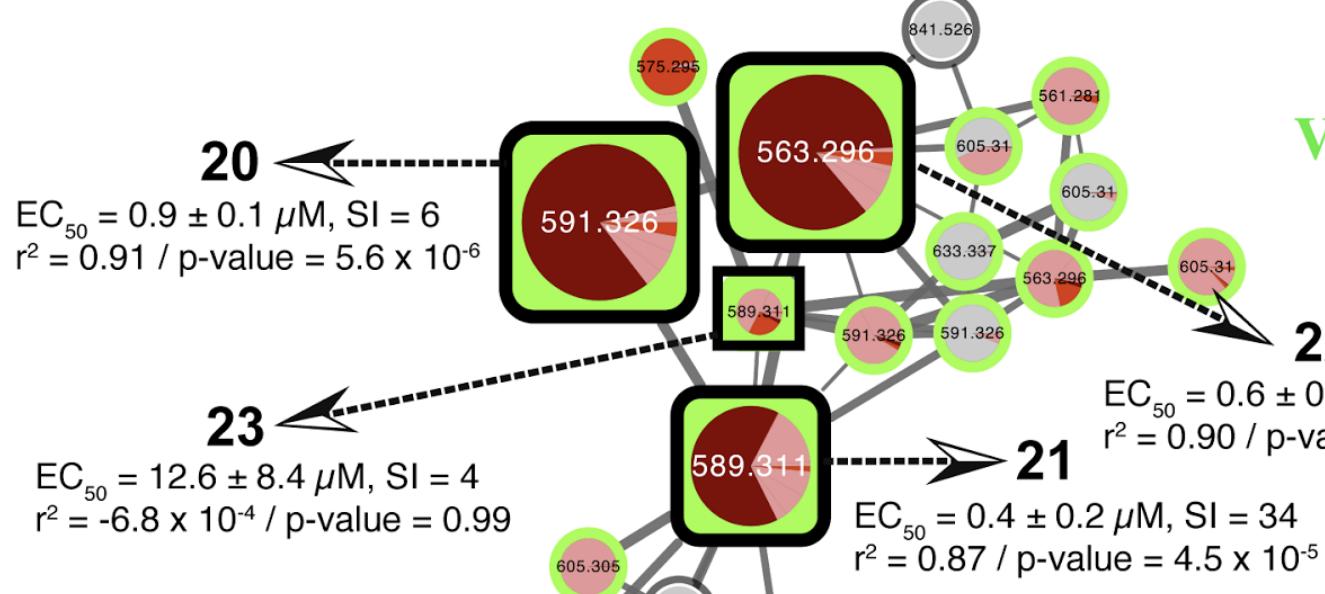
# MS based targeted isolation - NMR



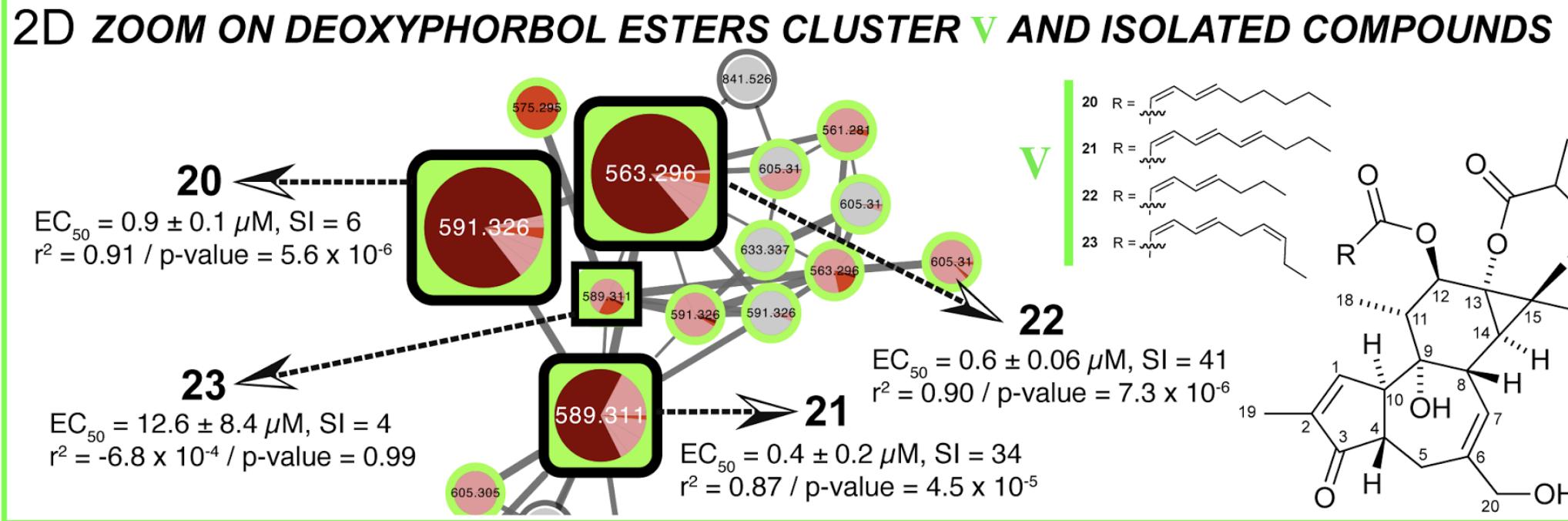
New family of deoxyphorbol esters

# Structures and anti-chikungunya virus activity

## DEOXYPHORBOL ESTERS CLUSTER



This discovery would not have been made if data and metadata was not public



# Example Data reuse 2.

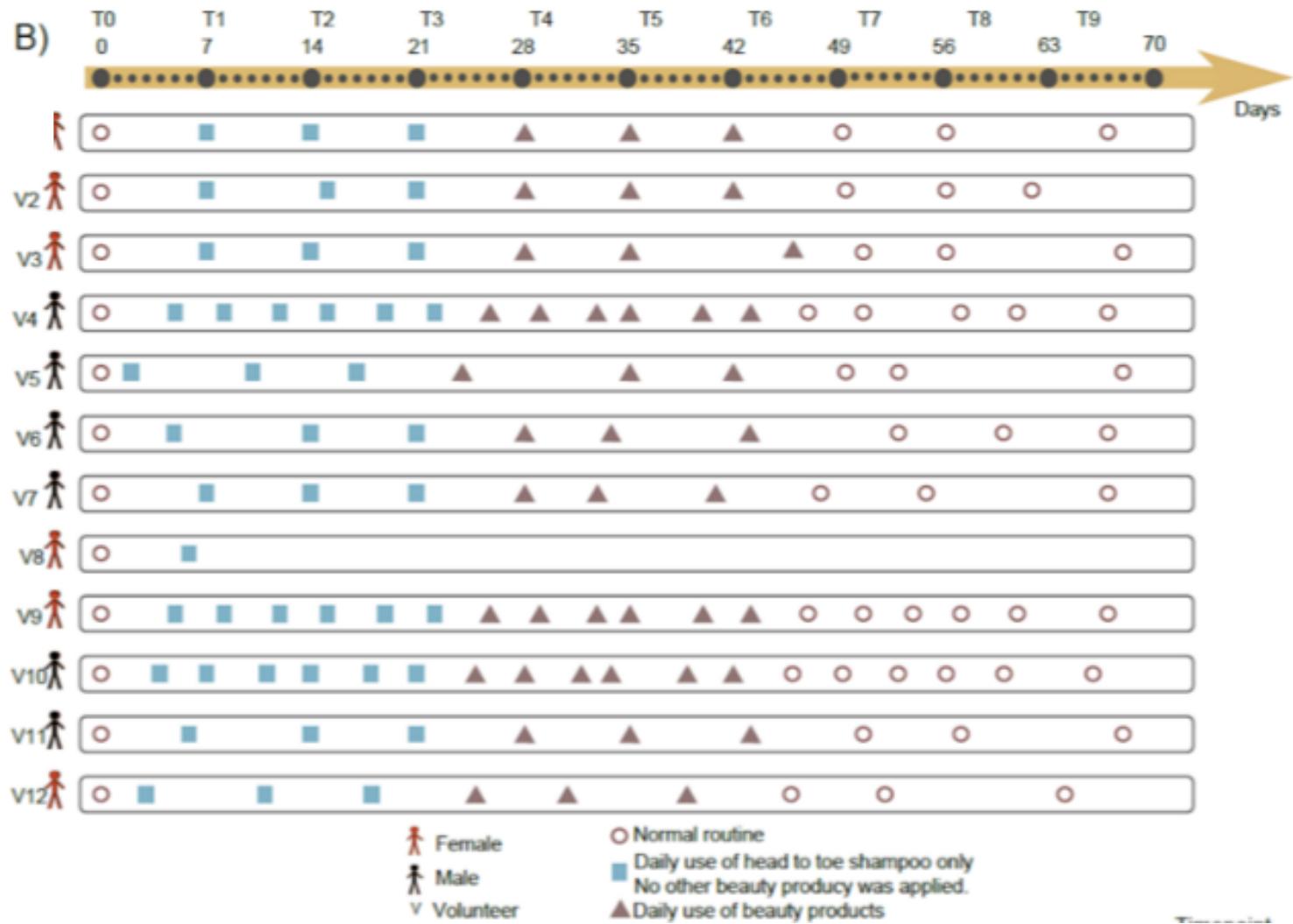
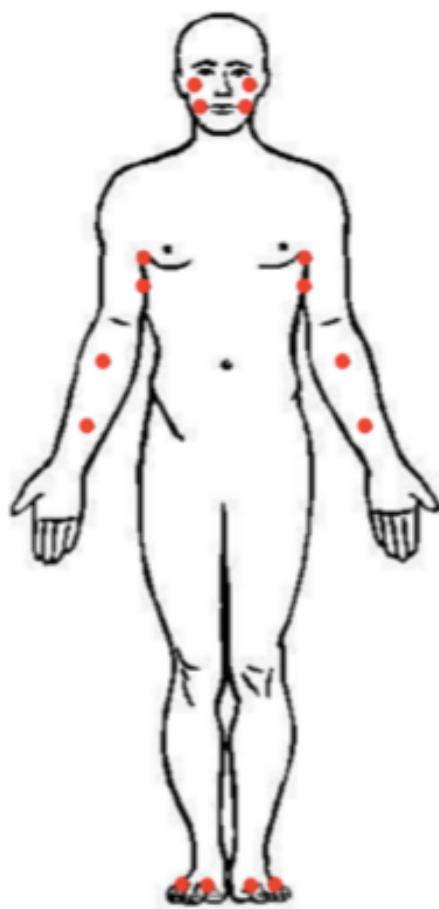
$m/z$  50

$m/z$  value

Red = high intensity  
Orange  
Yellow  
Green  
Light Blue  
Blue = low intensity

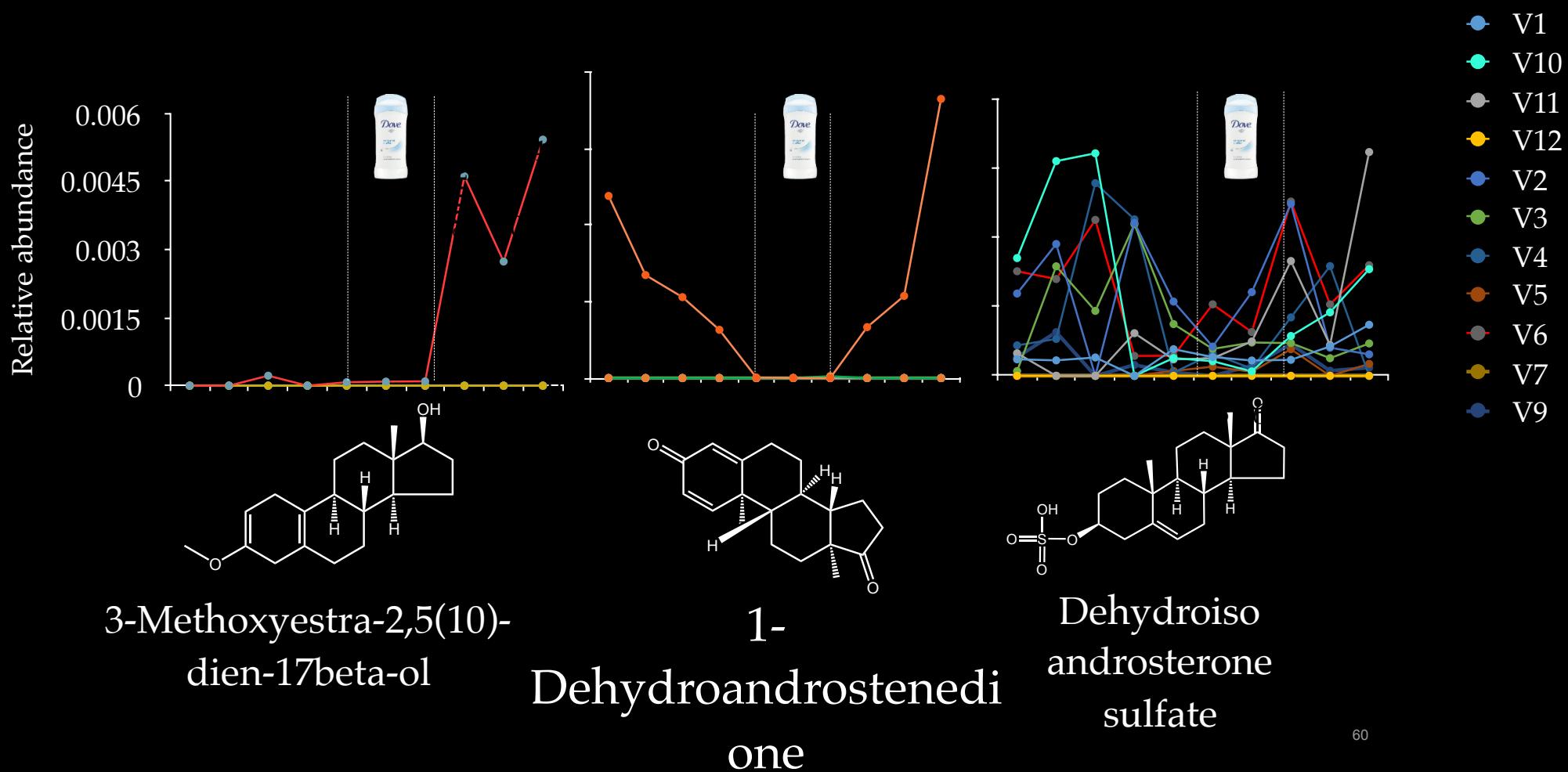


Does personal care impact the skin?



# Steroid molecular family in Armpits

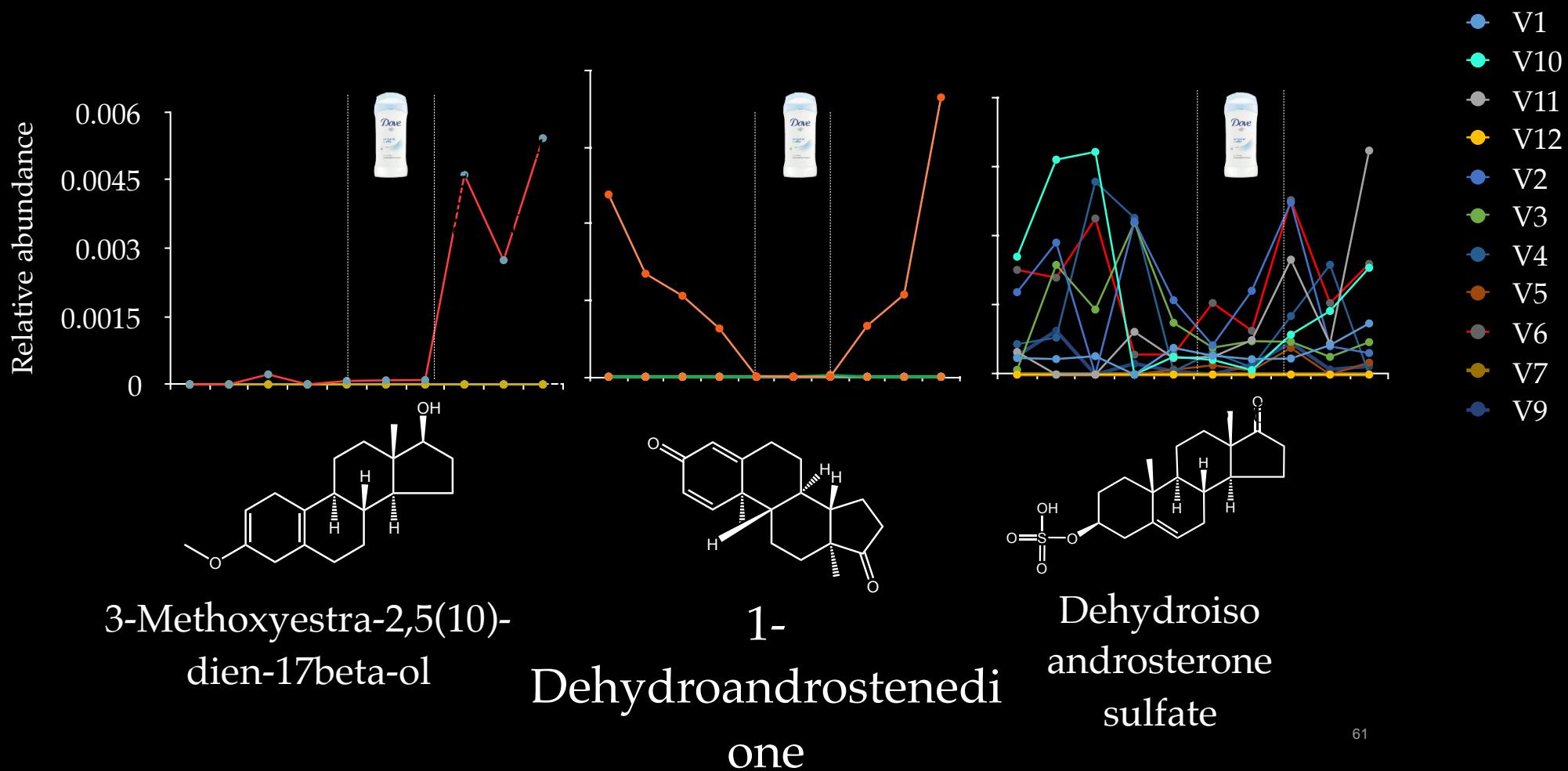
## Androsterone-response is person specific.



Influence of Androstenol and Androsterone on the Evaluation of Men of Varying Attractiveness Levels

# Steroid molecular family in Armpits

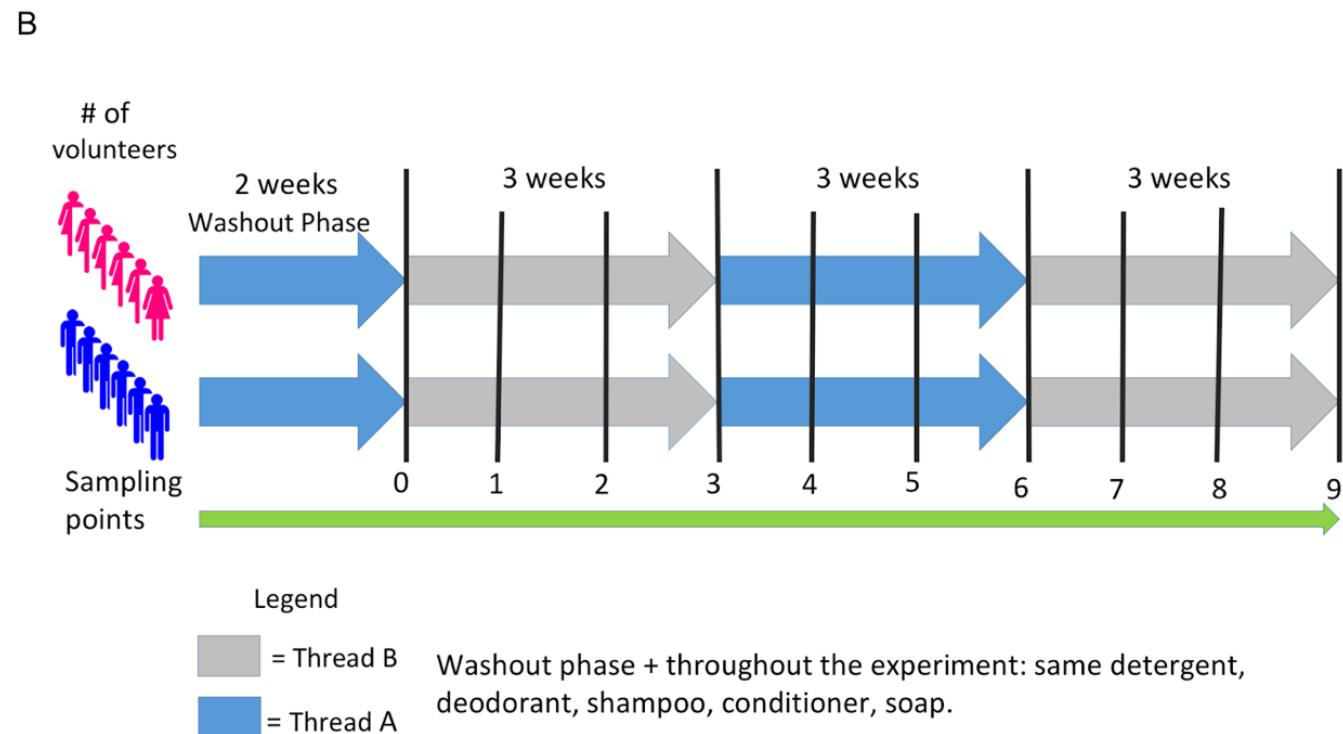
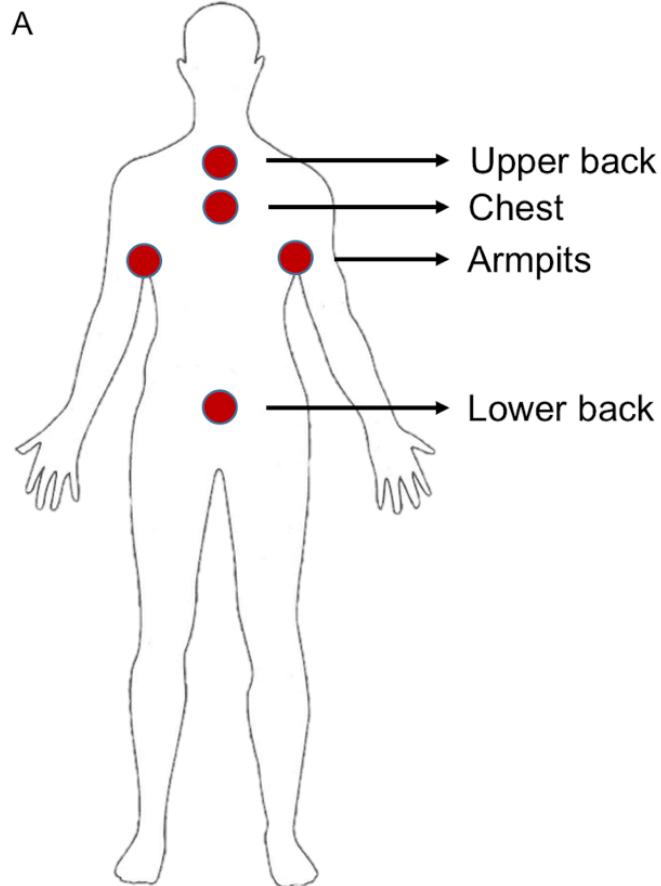
Had to get DEA approved, background checks etc. to get the standard so we could validate.



Influence of Androstenol and Androsterone on the Evaluation of Men of Varying Attractiveness Levels

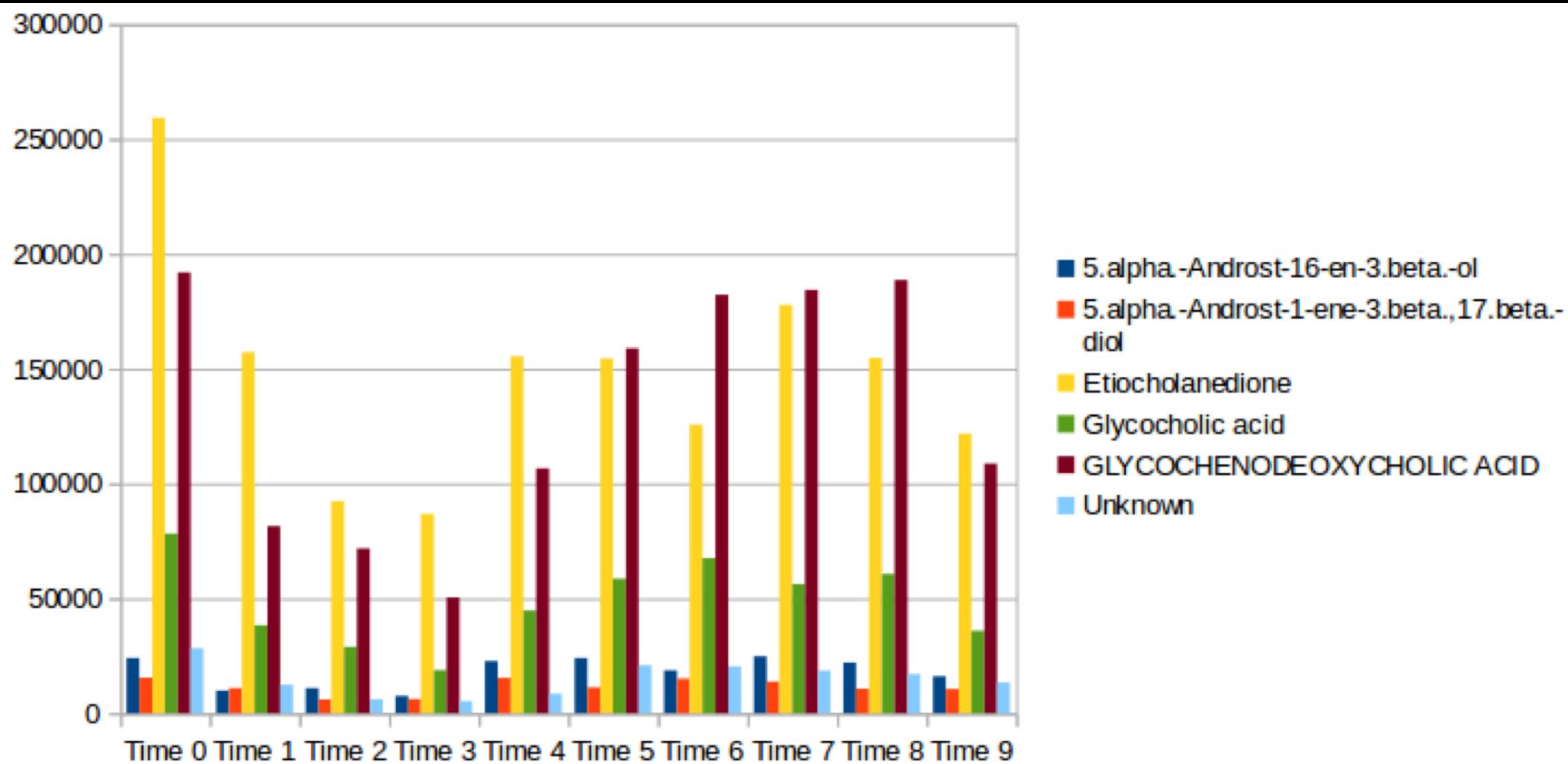
# Steroid molecular family in Armpits

Had to get DEA approved, background checks etc. to get the standard so we could validate.



# Steroid molecular family in Armpits

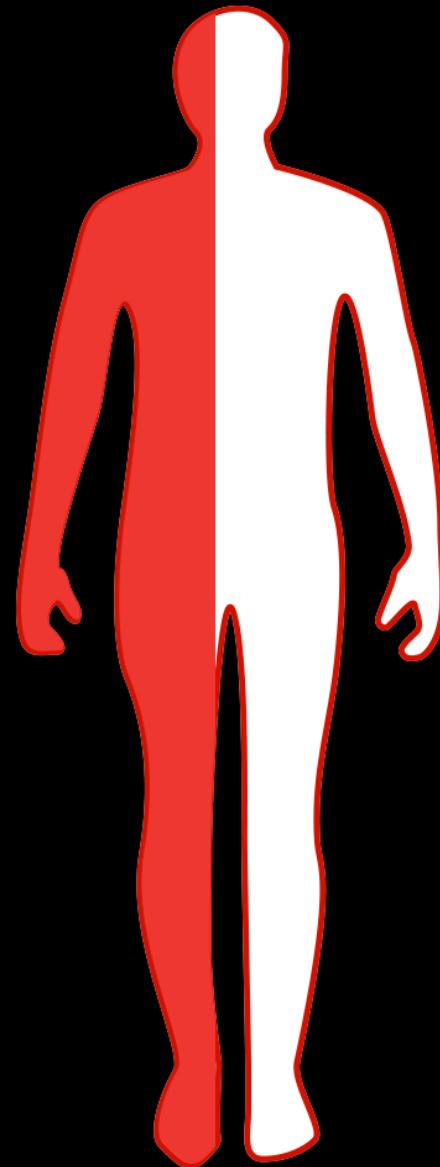
We could reuse data to validate.



# Example Data reuse 3.

**30 trillion  
human cells**

**43%**



**39 trillion  
microbial cells**

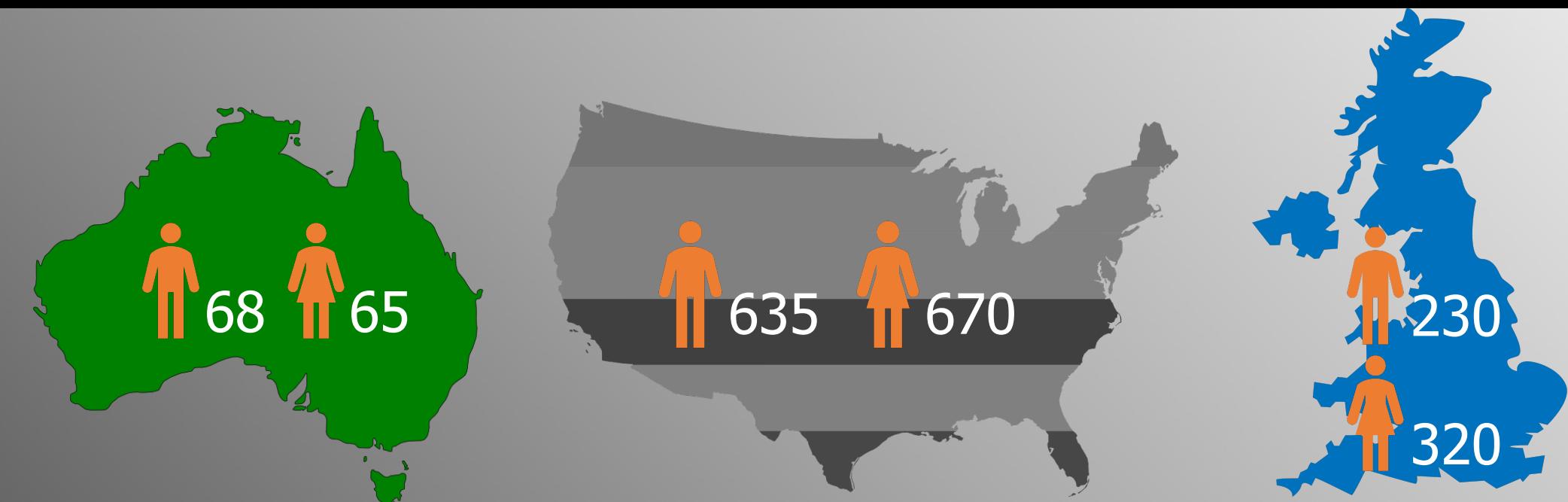


Citizen – scientist

Crowd funded - 15,000 sample collection

Paired microbiome & metabolome

Extensive metadata, self-reporting

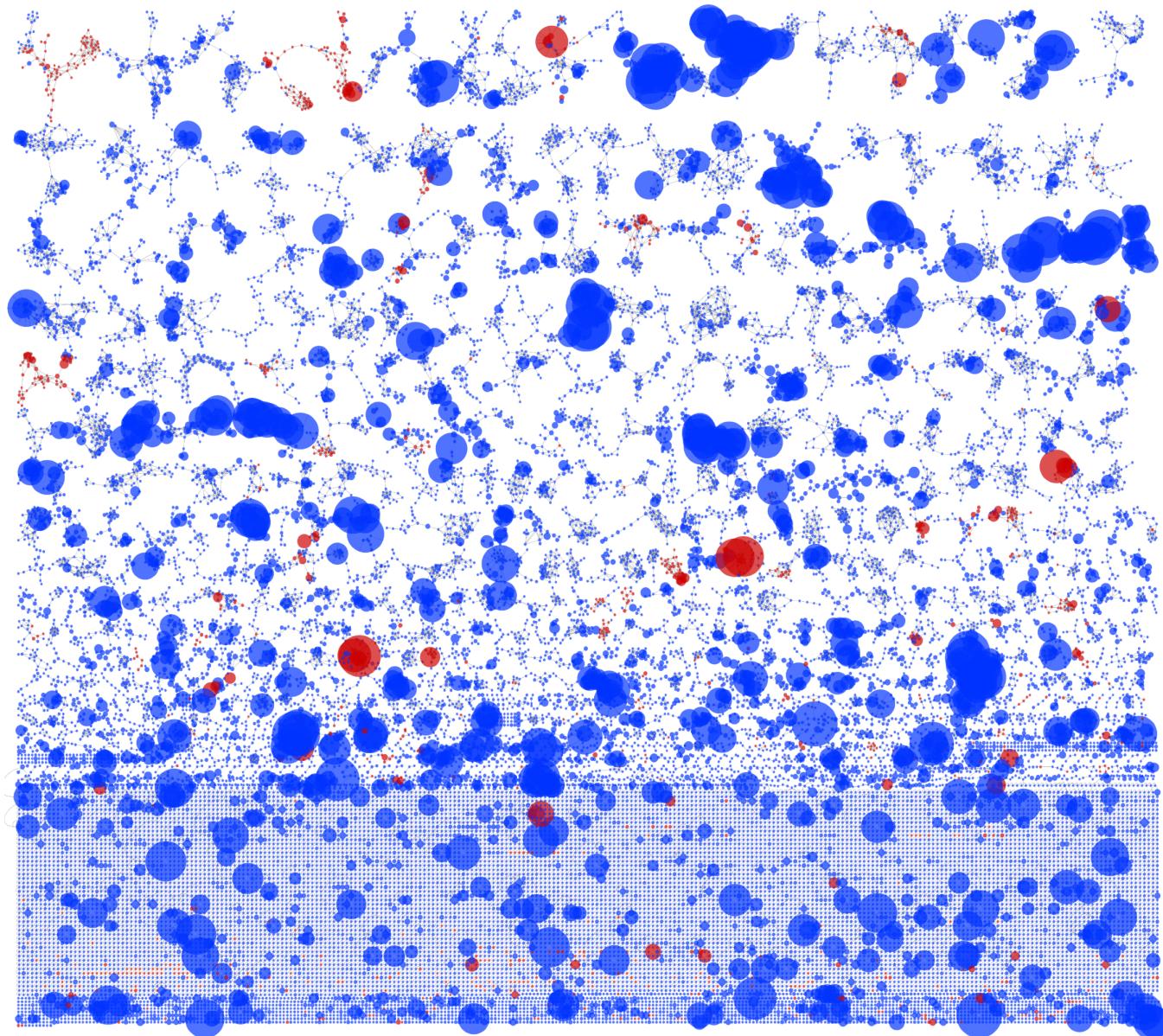


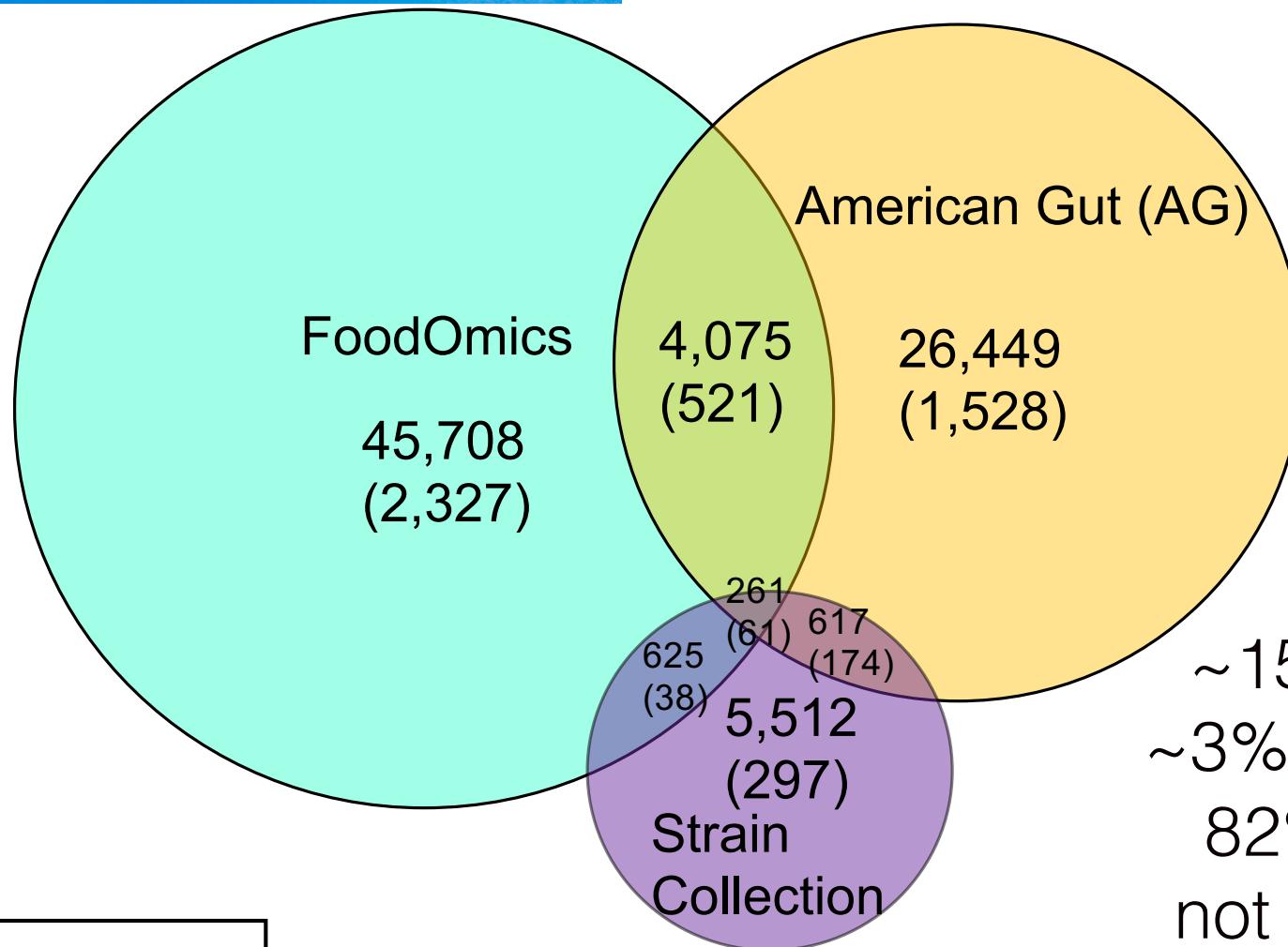
Do we see any microbial molecules?

# GNPS Networking of American Gut Samples

31,511 nodes

unannotated ← → more  
annotated frequent frequent

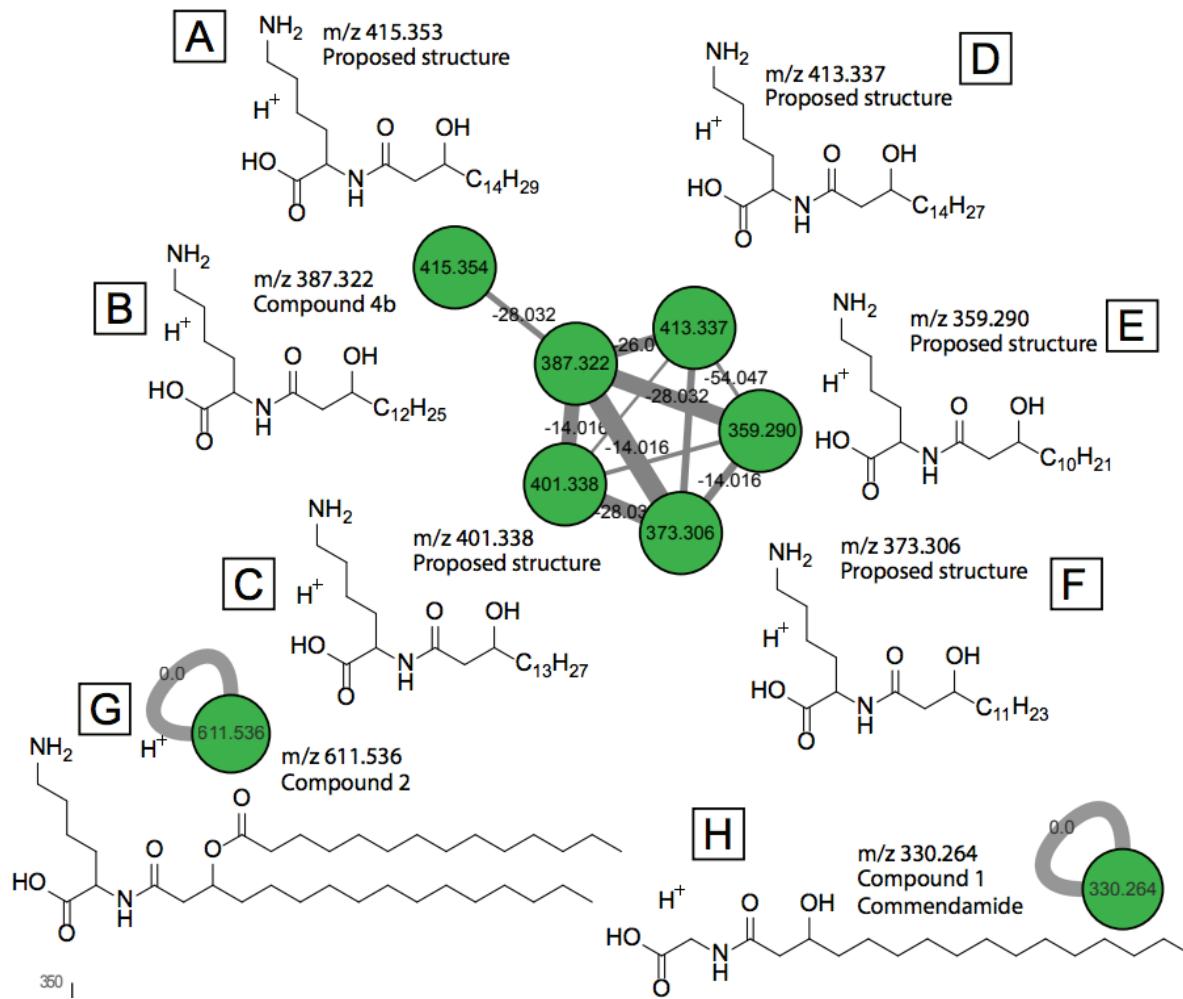




~15% Food  
~3% Microbes  
82% origin  
not matched

**NIH Human Microbiome Project**  
Microbial Reference Genomes



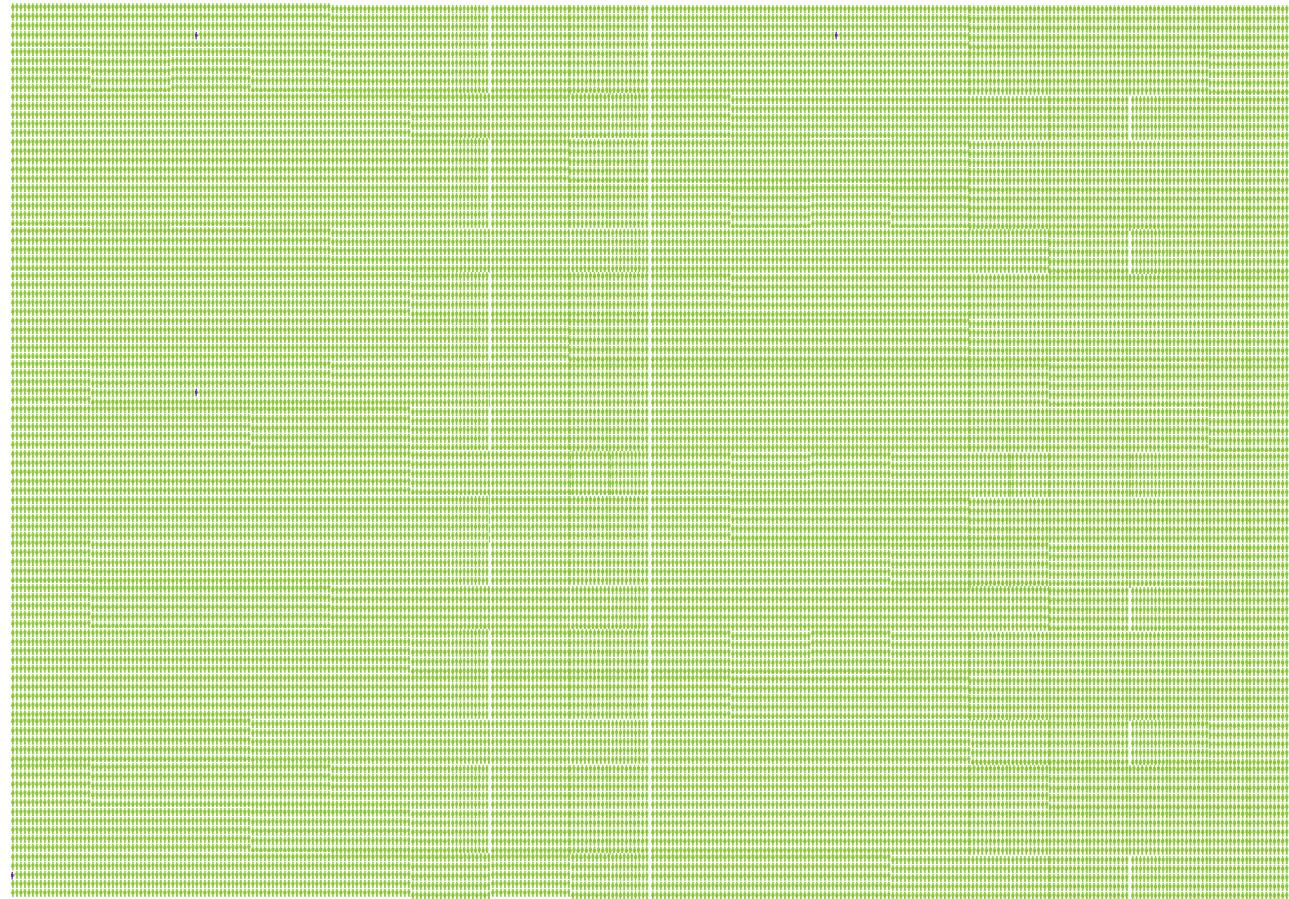
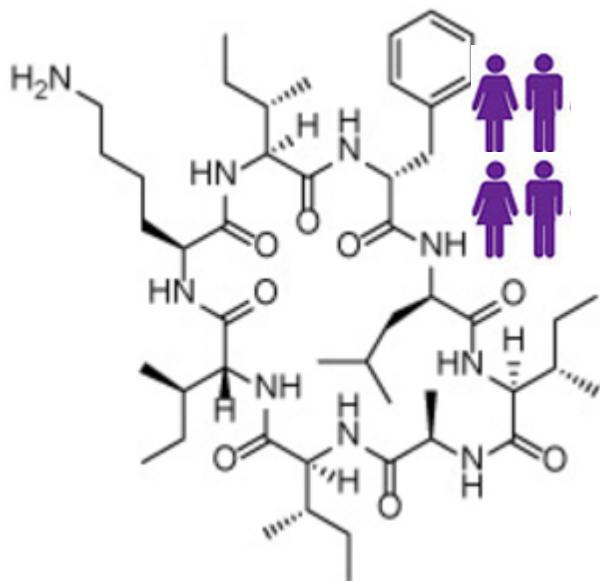


# Commensal bacteria make GPCR ligands that mimic human signalling molecules

Louis J. Cohen, Daria Esterhazy, Seong-Hwan Kim, Christophe Lemetre, Rhiannon R. Aguilar, Emma A. Gordon, Amanda J. Pickard, Justin R. Cross, Ana B. Emiliano, Sun M. Han, John Chu, Xavier Vilafarres, Jeremy Kaplitt, Aneta Rogoz, Paula Y. Calle, Craig Hunter, J. Kipchirchir Bitok & Sean F. Brady

regulates glucose homeostasis

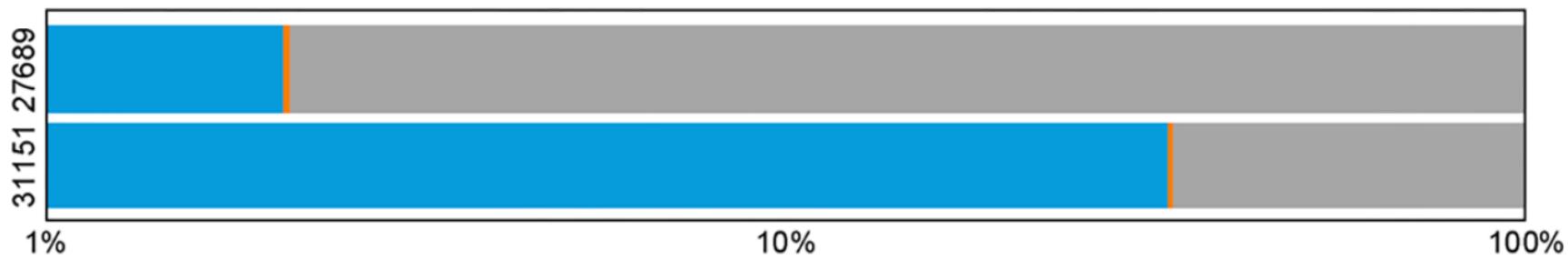
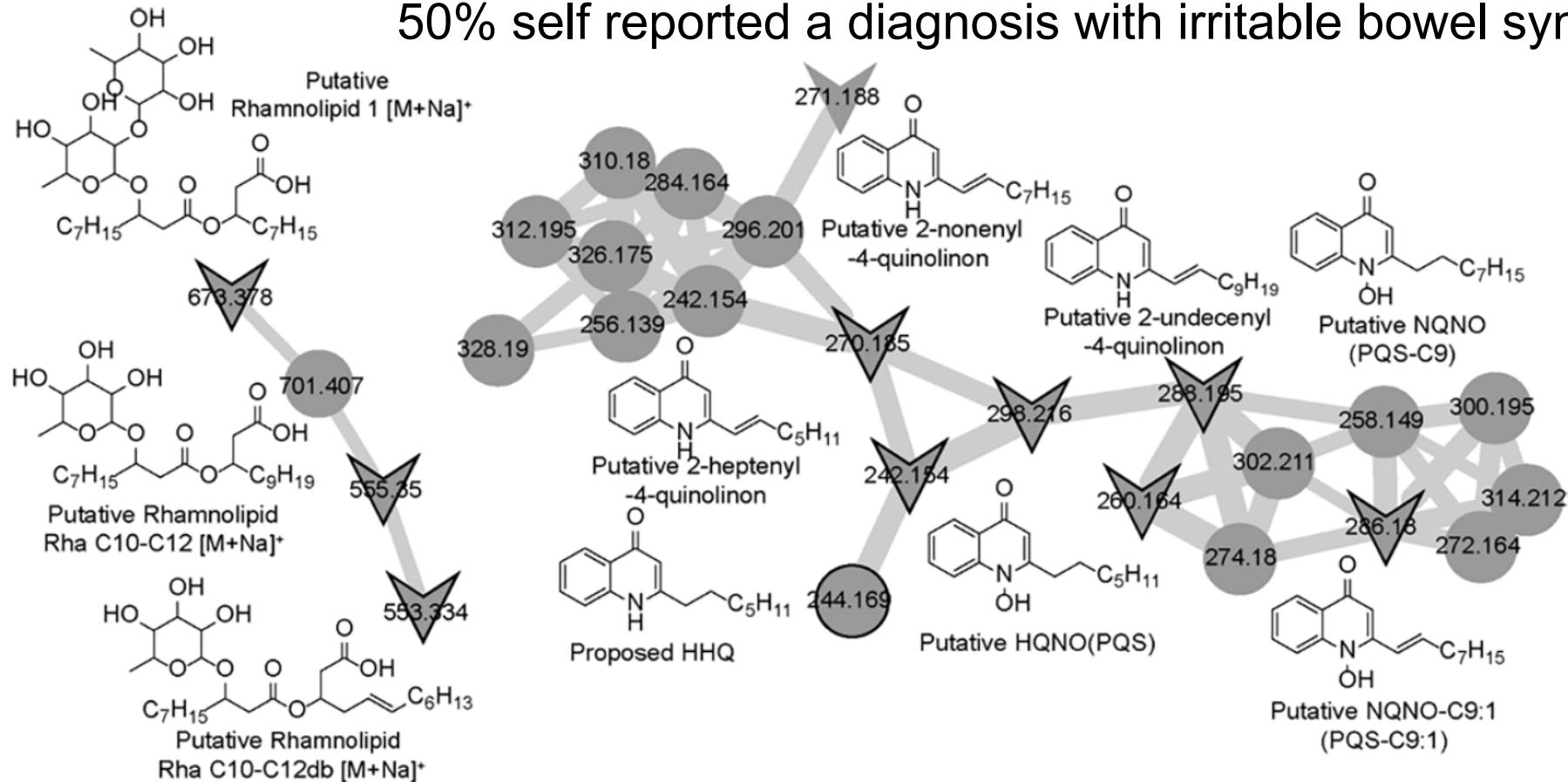
# Matches to microbial molecules annotated by the GNPS community



Matched to 4 out of the 2100 samples

# Matches to microbial molecules annotated by the GNPS community

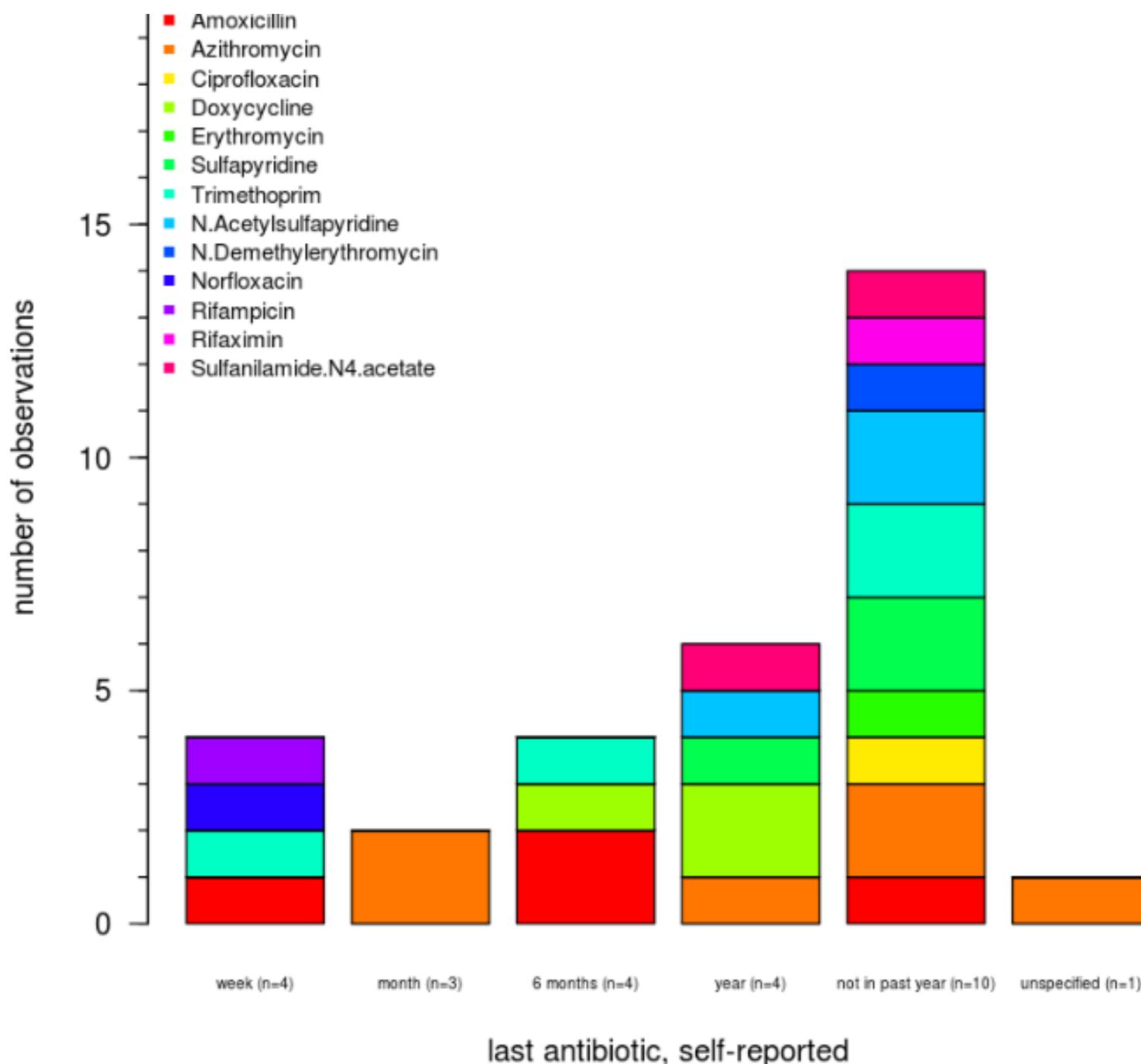
*Pseudomonas aeruginosa*-associated molecules.  
50% self reported a diagnosis with irritable bowel syndrome.



■ f\_Pseudomonadaceae   ■ f\_Pseudomonadaceae;g\_Pseudomonas   ■ Other taxa

Melnik

>80% of people where we detected antibiotics had not taken antibiotics in >6 months.

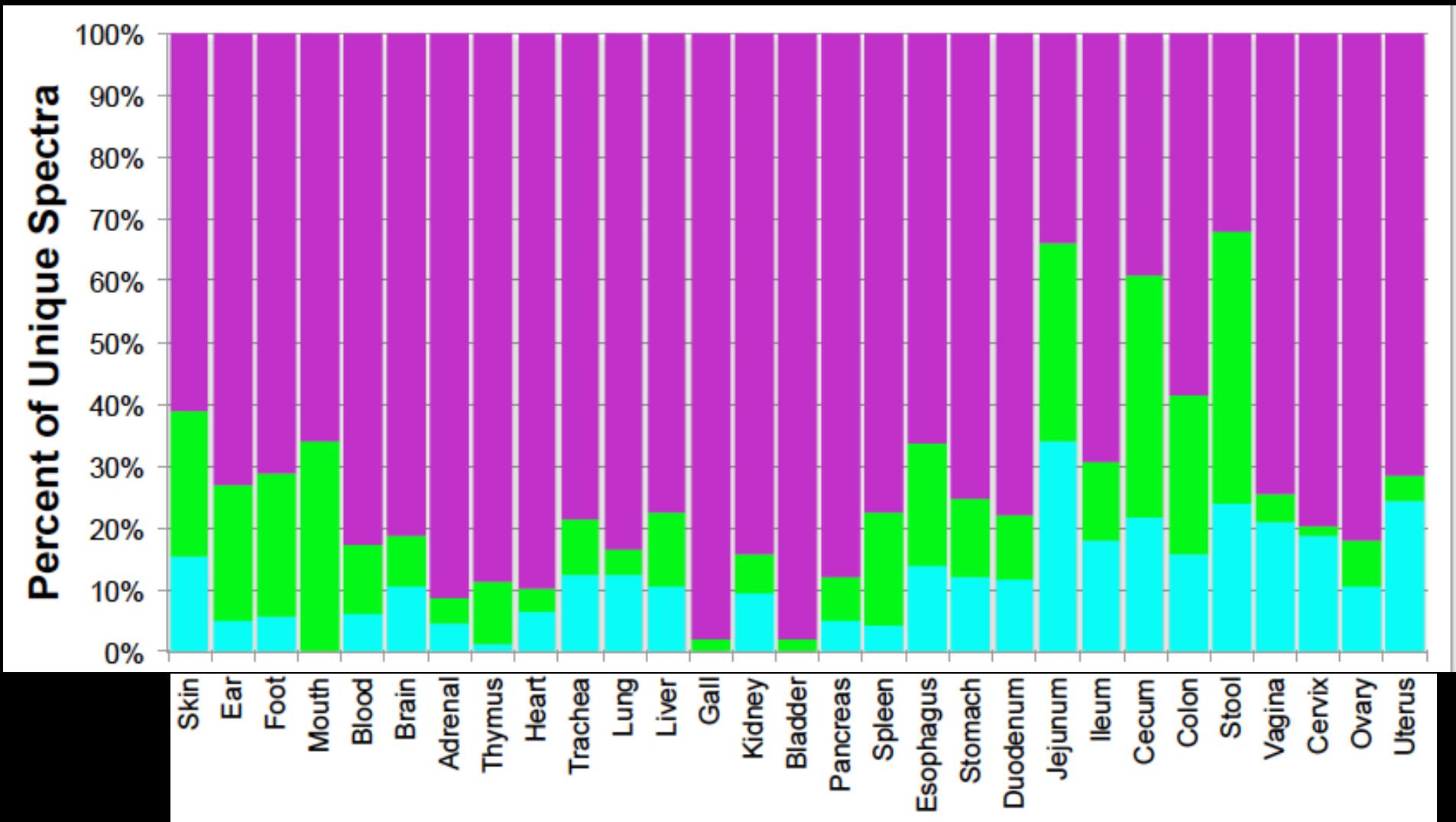


Example Data reuse 4.

We are not mono-colonized, we colonized by 1000s of microbes.

What impact does a “normal” microbiome have?

As much as 70% of the chemistry is different due to microbes  
What are these microbial molecules?



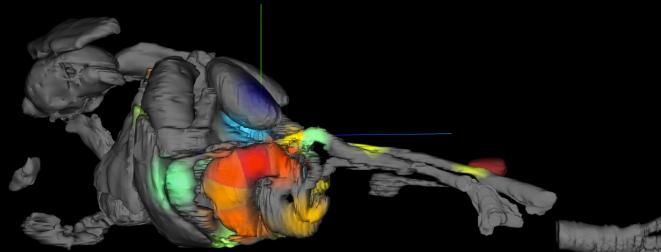
-Microbes

+Microbes

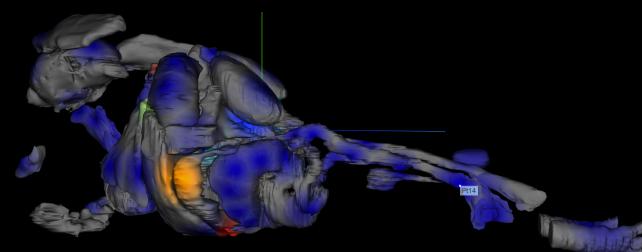
Shared

Paper will be submitted in a week or so

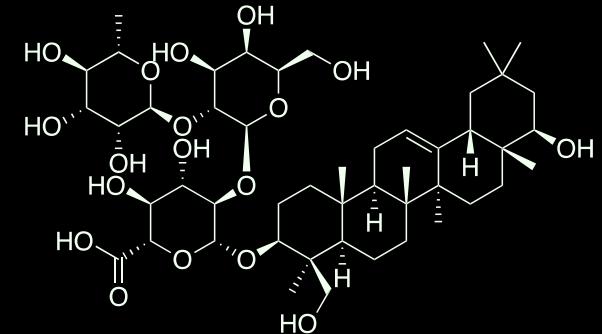
# Microbiome-Mediated Metabolism of Food



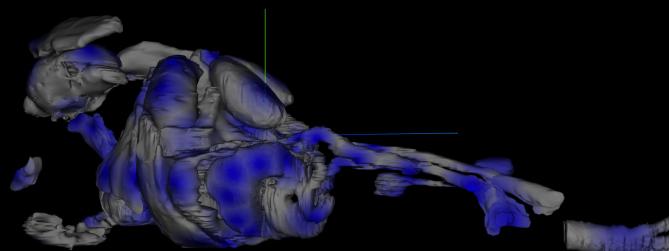
-Microbes



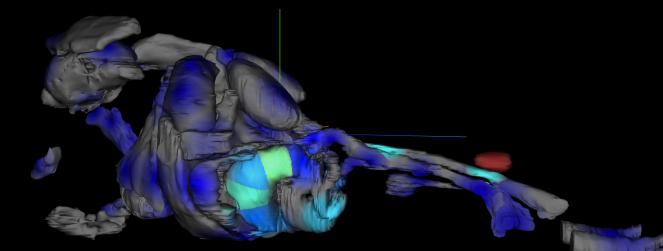
+Microbes



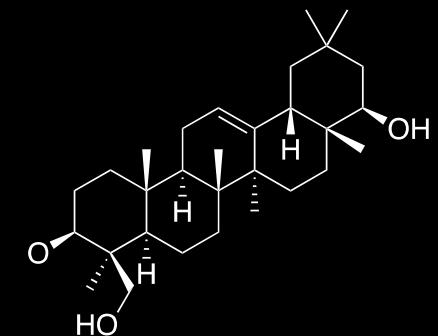
Soyasaponin



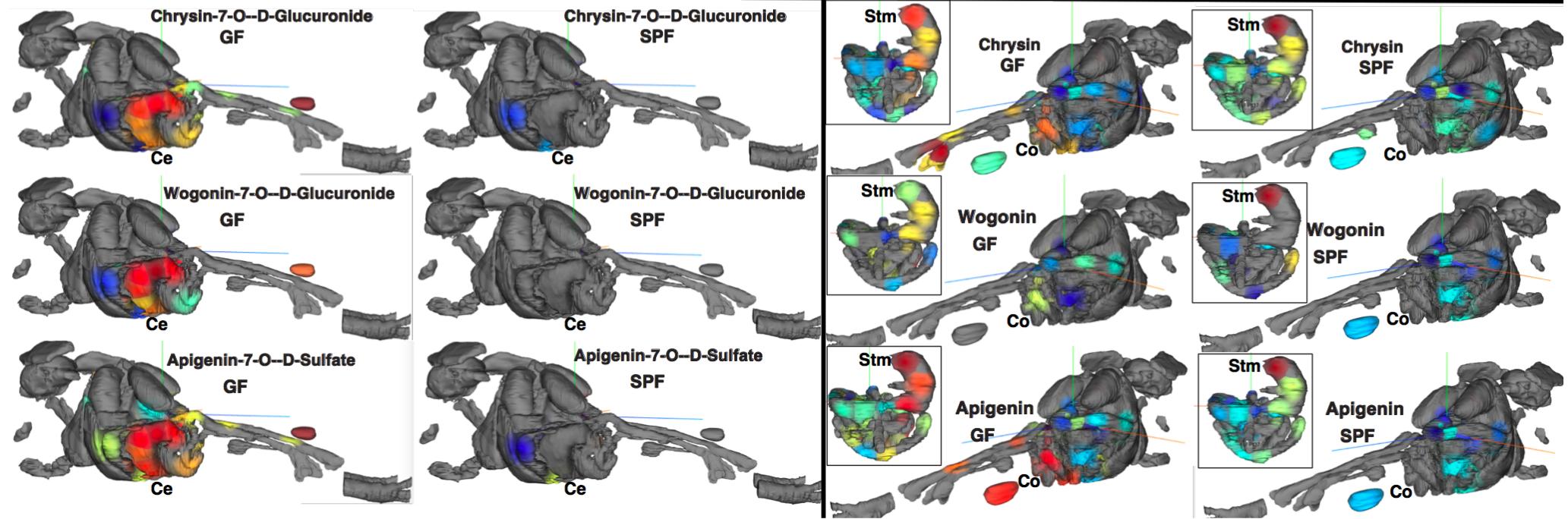
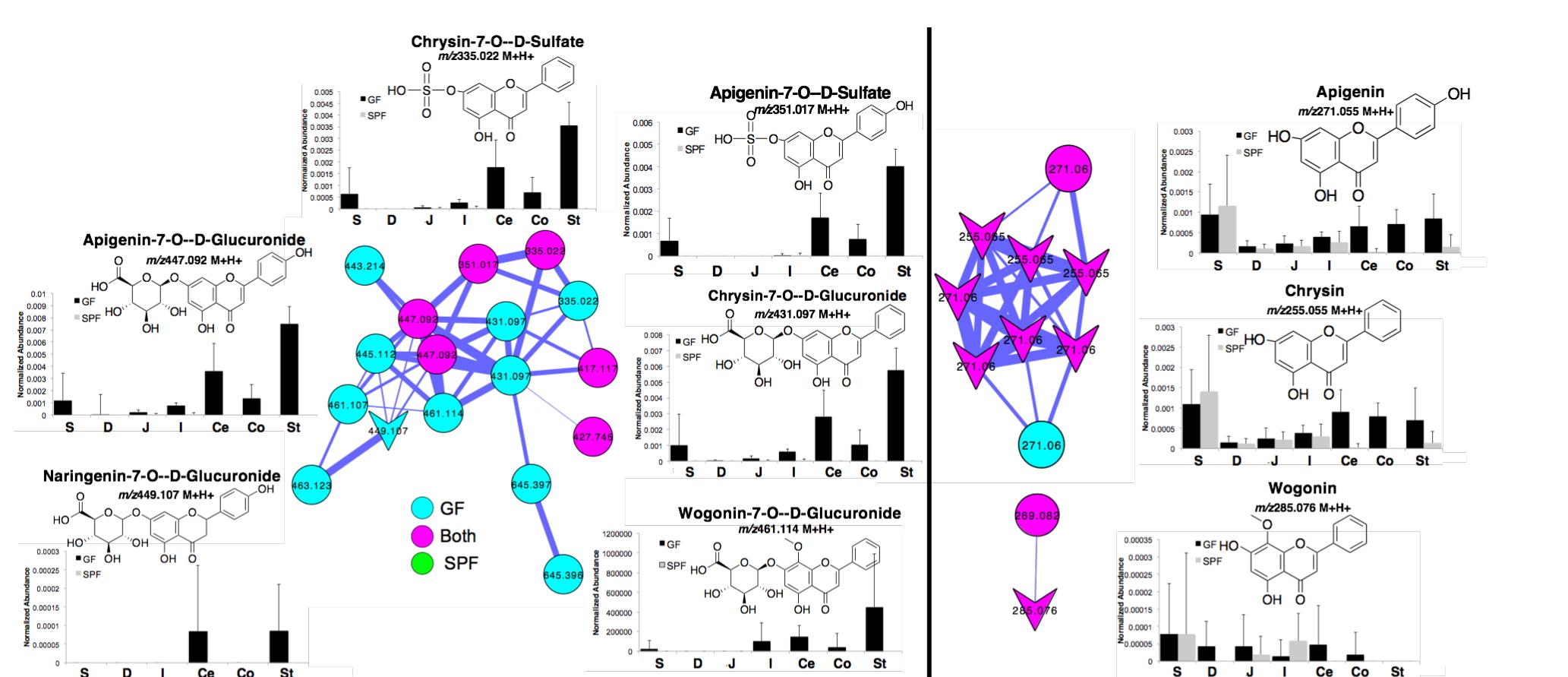
-Microbes

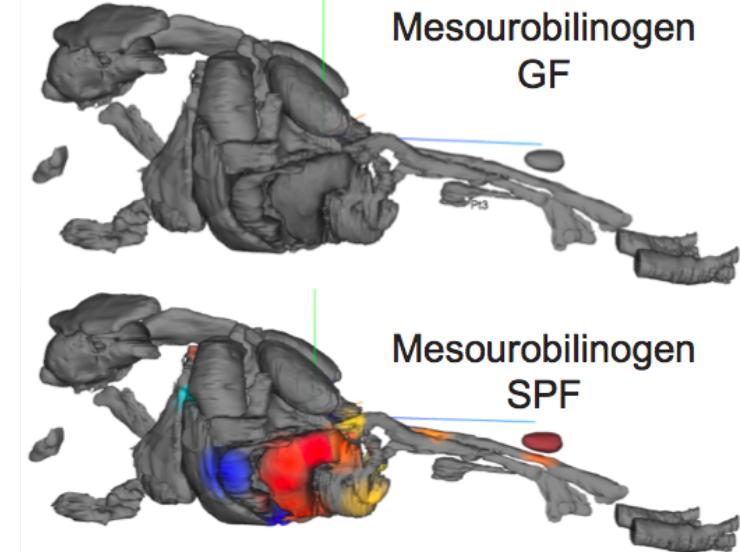
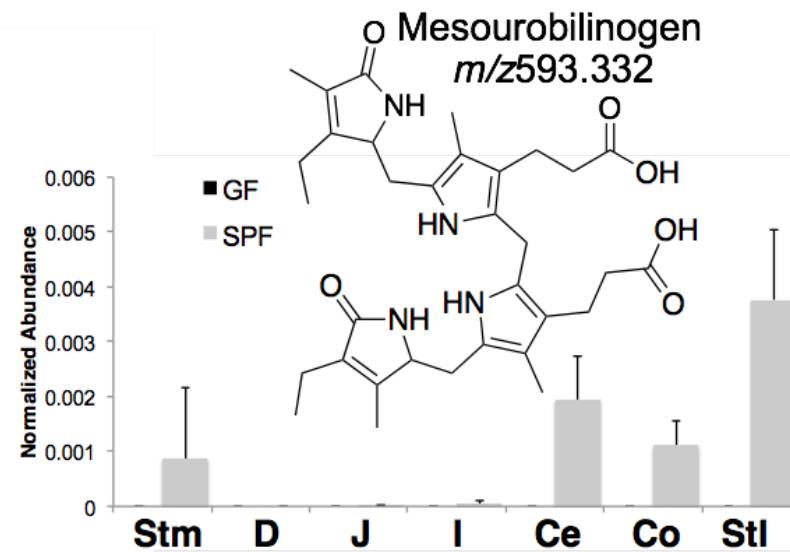
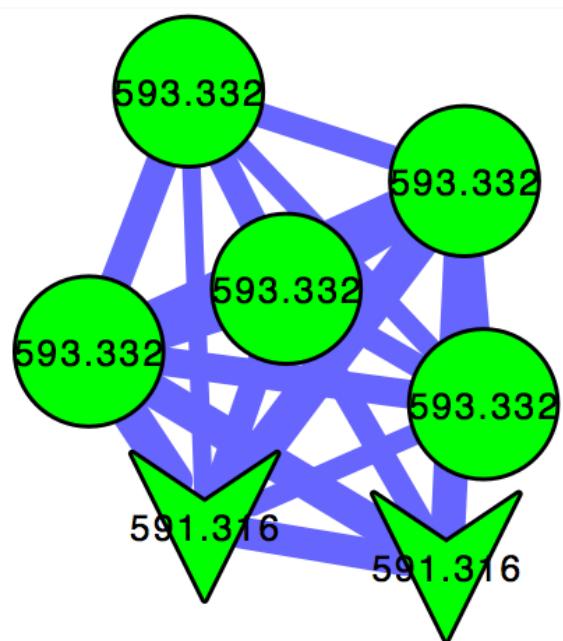


+Microbes

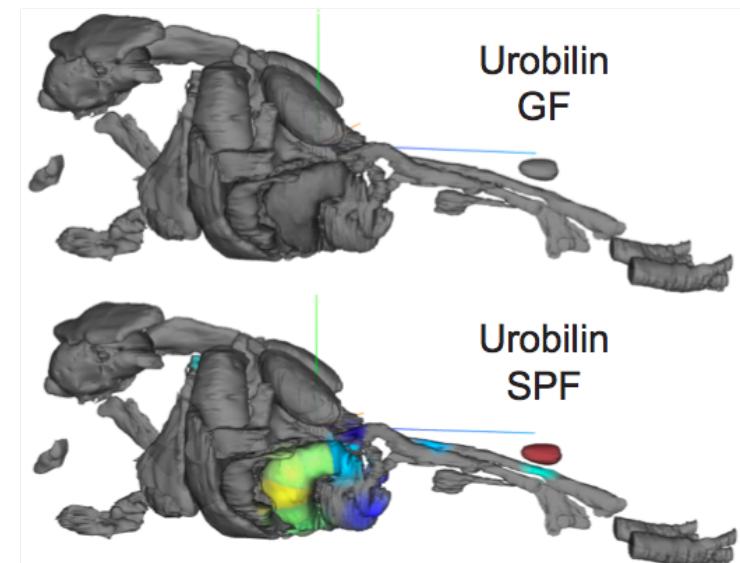
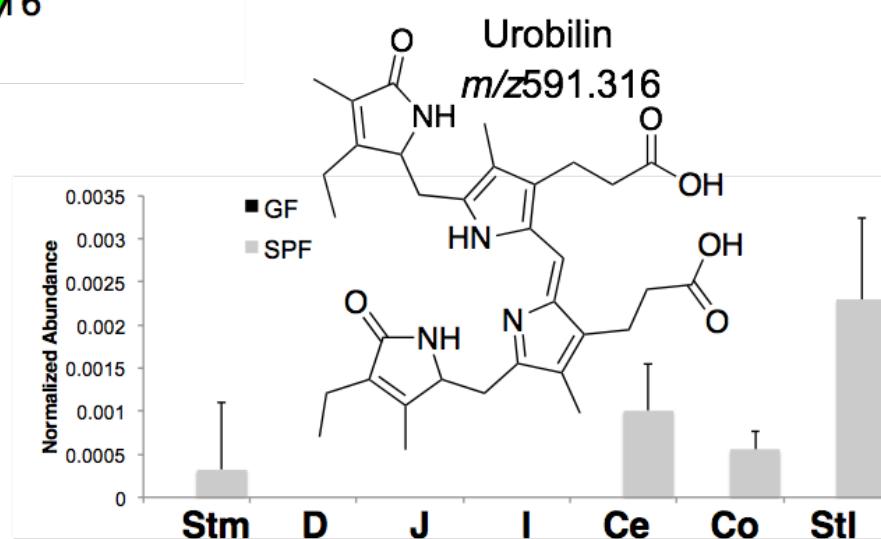


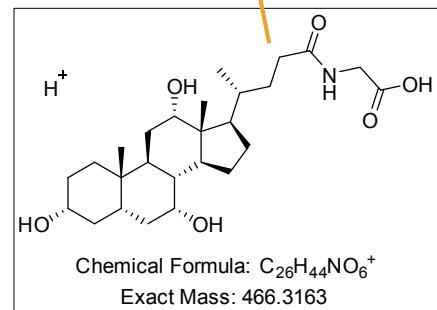
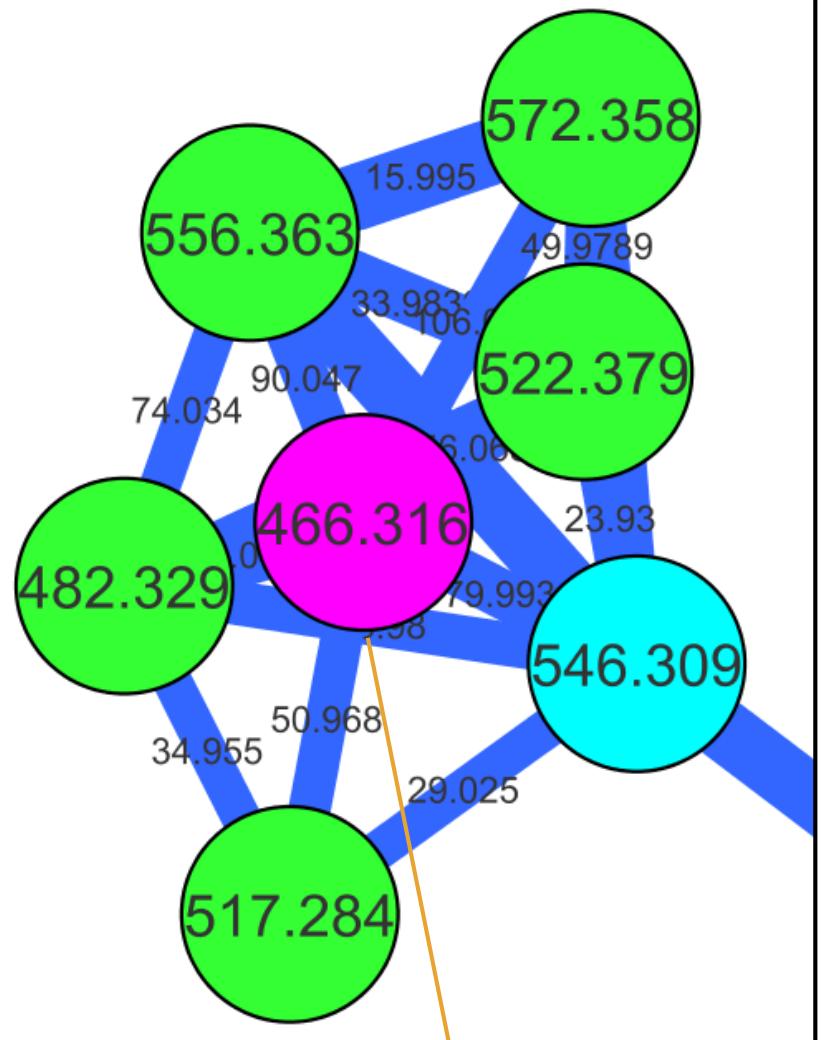
Soyaspongenol





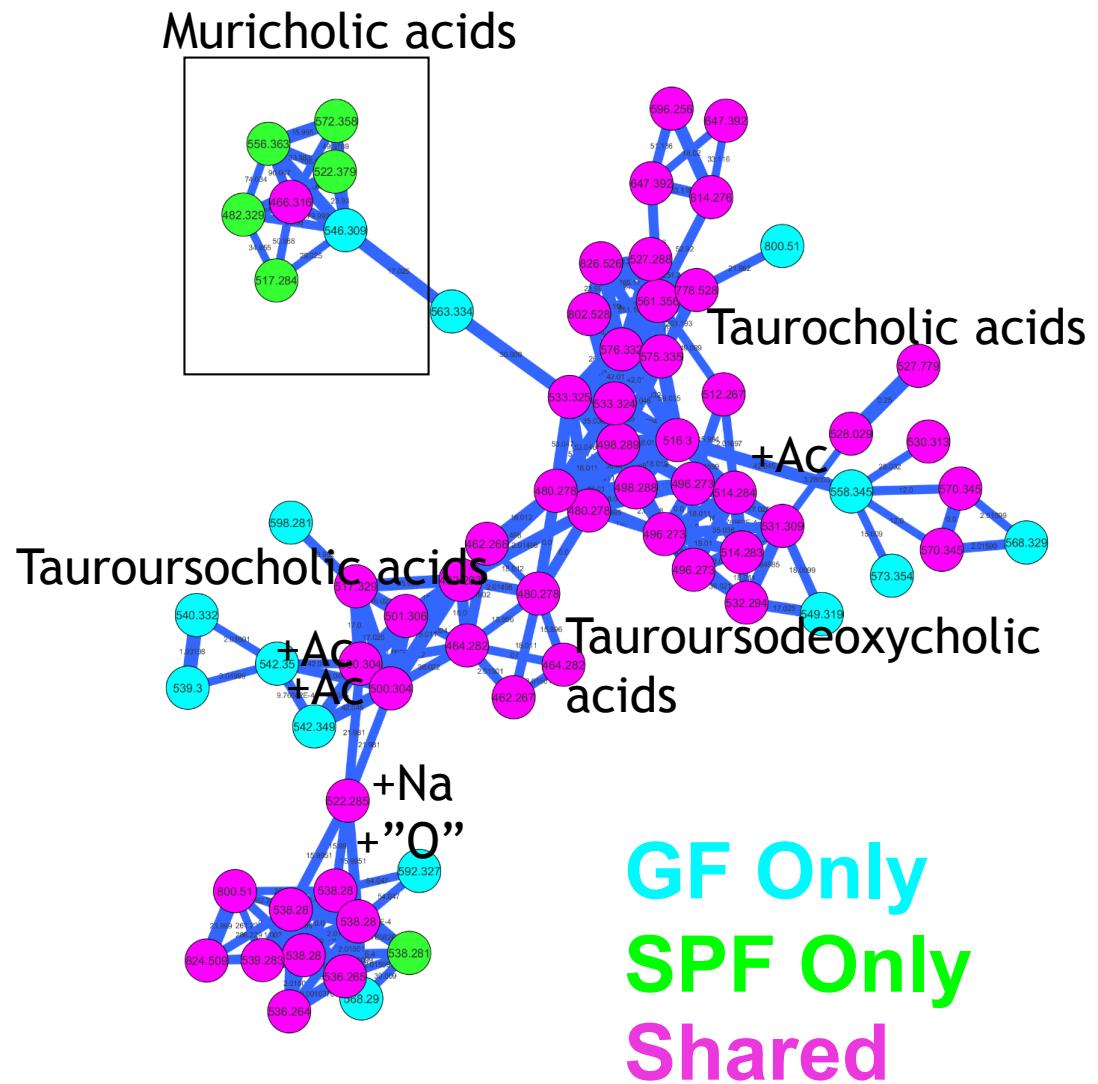
- GF
- Both
- SPF





Muriglycocholic acid

# Duodenum

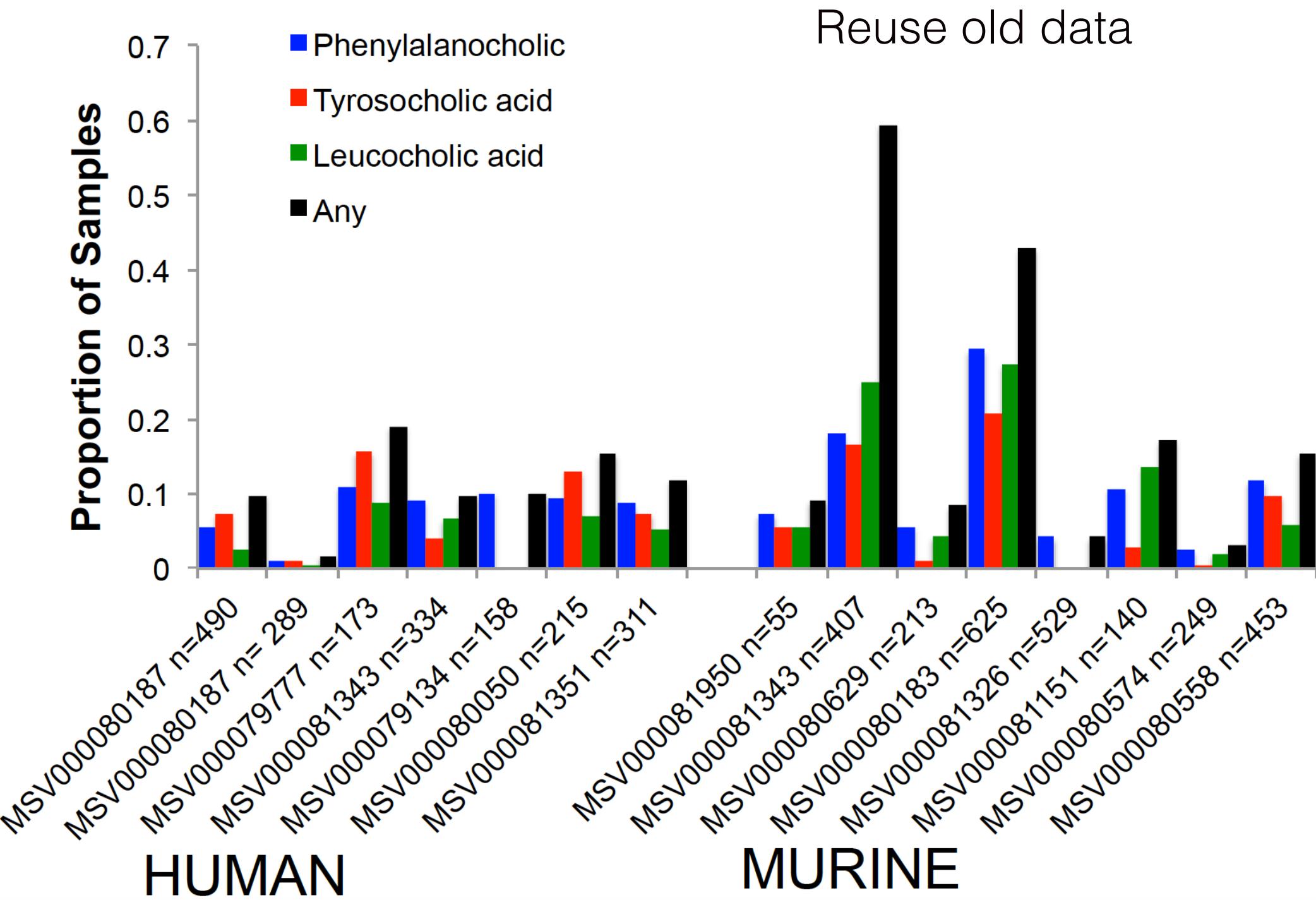


# Are they found in humans?

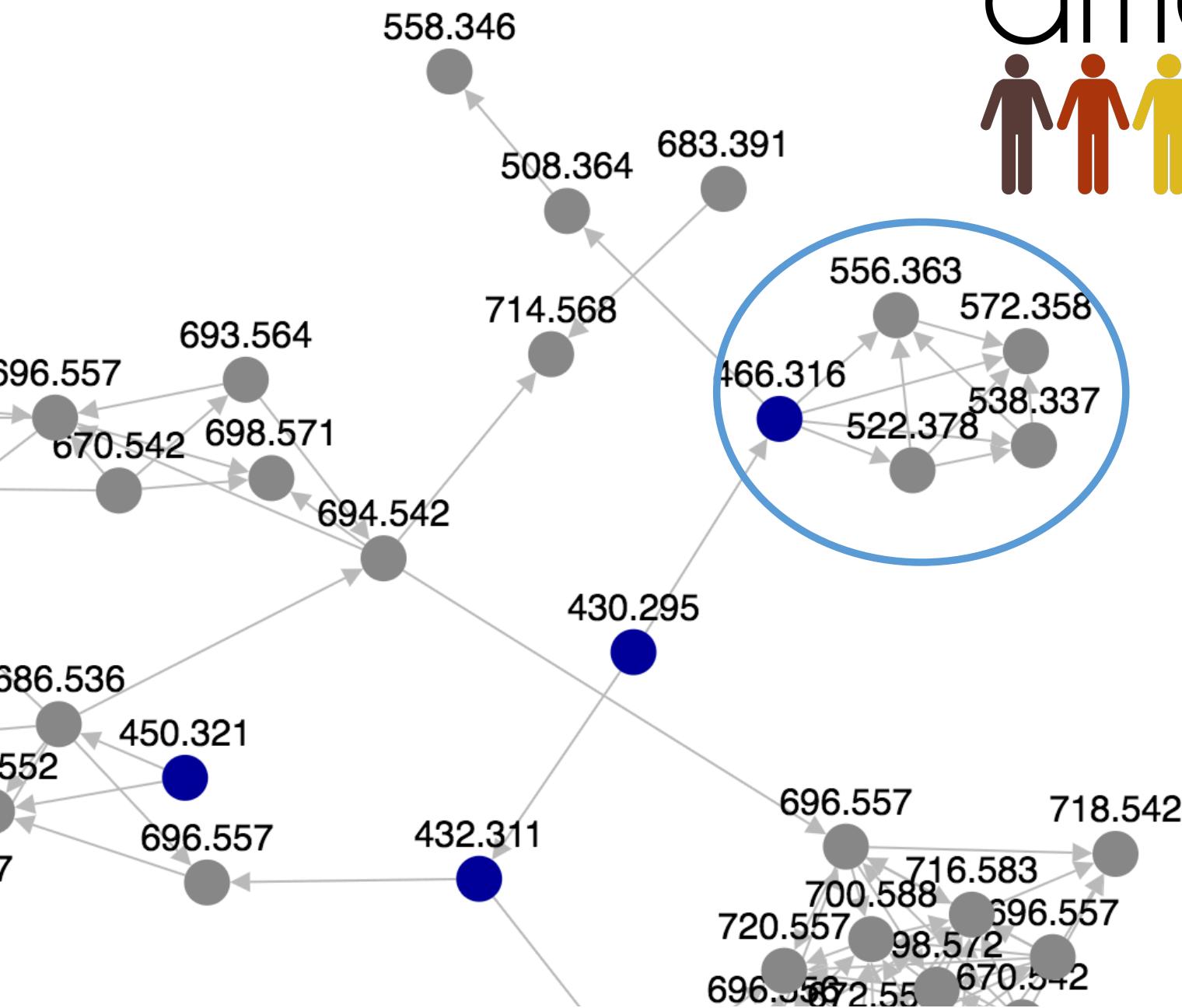
The screenshot shows the NCBI BLAST search interface. At the top, there are links for NIH, U.S. National Library of Medicine, NCBI National Center for Biotechnology Information, and a sign-in option. Below the header, it says "BLAST® > blastp suite". The main search area has tabs for blastn, blastp, blastx, tblastn, and tblastx. A large search input field is labeled "Enter Query Sequence" and "Enter accession number(s), gi(s), or FASTA sequence(s)". Below the input field are options for "Or, upload file", "Job Title", and "Enter a descriptive title for your BLAST search". There is also a checkbox for "Align two or more sequences". A prominent feature is a large black rectangular box containing the text "Search Single Spectrum Over Datasets Available [here](#)". To the right of the search form, there is a sequence of DNA bases: GACAGACATGACTTGGATTCCCCAGGAAGGAGTTGGCAACCCAGTCAAAGGC... and a link to "Reset page". On the far right, there is a "Where is this seq." section with a question mark and a "Help" link. At the bottom left, there is a "Choose Search Set" section with dropdowns for "Database" (set to "Non-redundant protein sequences (nr)"), "Organism" (optional), "Exclude" (checkboxes for "Models (XM/XP)" and "Uncultured/environmental sample sequences"), and "Entrez Query" (optional). A magnifying glass icon is overlaid on the search interface, and three cylindrical databases labeled "Database", "Knowledgebase", and "Retrieval infrastructure" are shown at the bottom.



Retrieval infrastructure

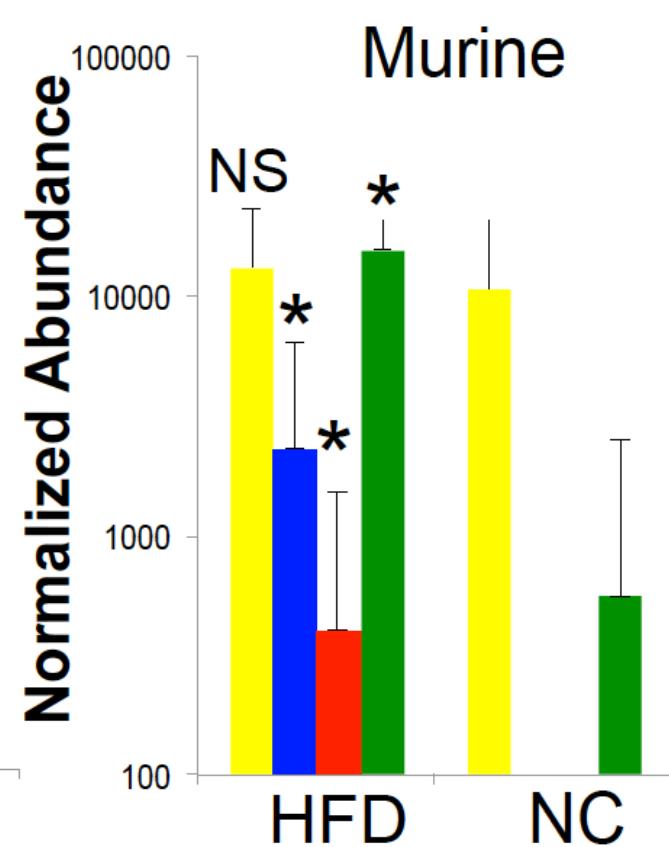
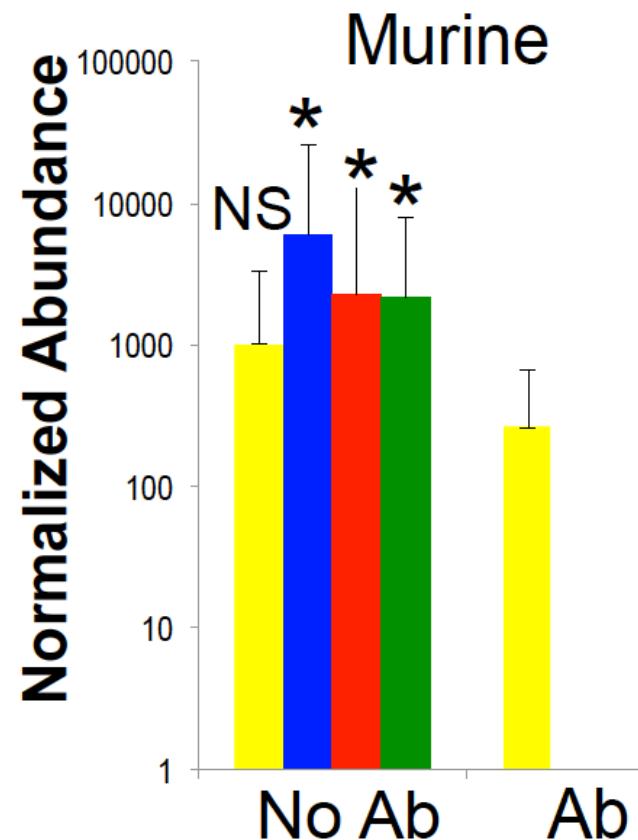
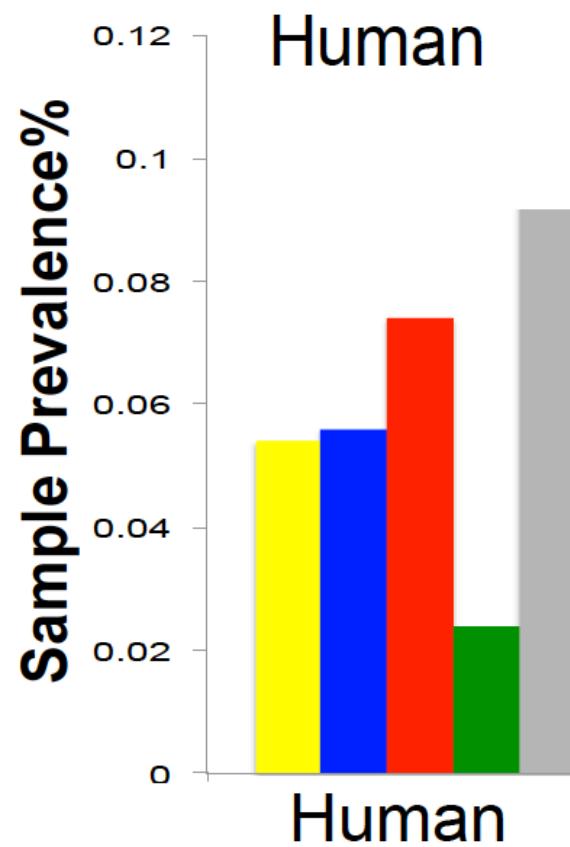


Reuse old data: Are they found in humans?



# Looking at the metadata of published data

■ Glycocholic acid ■ Phenylalanocholic acid  
■ Tyrosocholic acid ■ Leucocholic acid ■ All

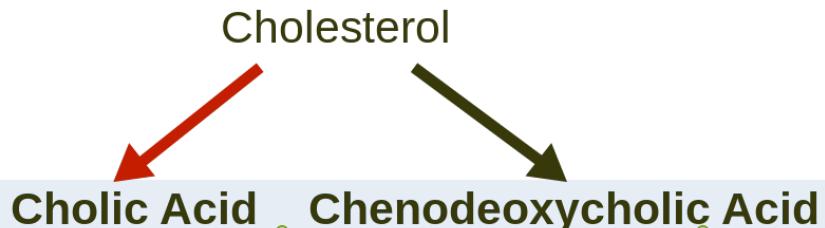


Generally reduced with antibiotic use, diarrhea increased, in high fat diet increased, viral infection increased

Bacteria are responsible for

- 1) Dehydroxylation
- 2) Deconjugation
- 3) **Conjugation**

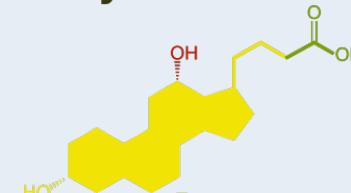
## Primary Bile Acids



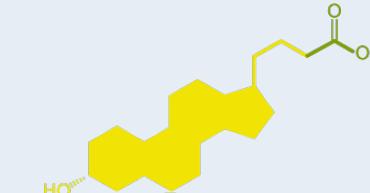
dehydroxylation  
in gut by bacteria

## Secondary Bile Acids

### Deoxycholic Acid

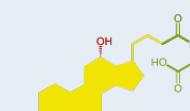


### Lithocholic Acid



conjugation  
in liver

### Glyco-cholic Acid



### Taur-cholic Acid



### Glycodeoxy-cholic Acid



### Taurodeoxy-cholic Acid



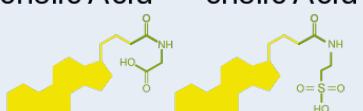
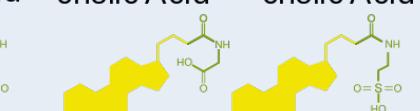
### Glycolitho-cholic Acid



### Taurolitho-cholic Acid



### Glycochenodeoxy-cholic Acid



## Conjugated Bile Acids

Chemical Formula: C<sub>33</sub>H<sub>50</sub>NO<sub>6</sub><sup>+</sup>

Exact Mass: 556.3633

Chemical Formula: C<sub>33</sub>H<sub>50</sub>NO<sub>6</sub><sup>+</sup>

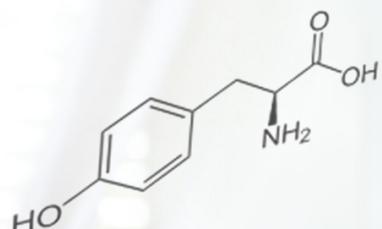
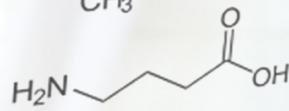
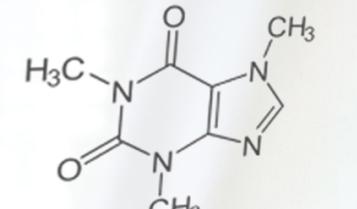
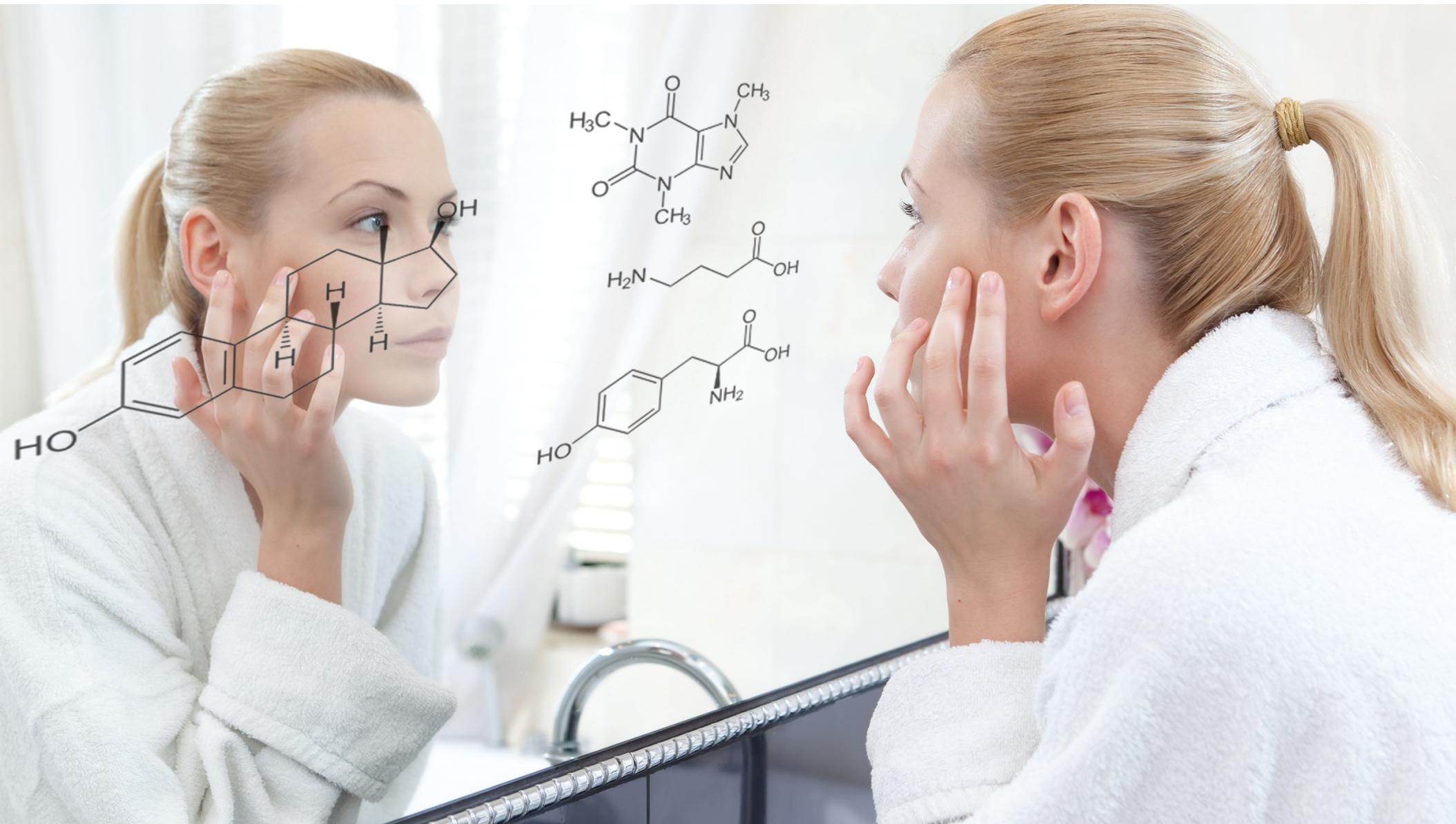
Chemical Formula: C<sub>33</sub>H<sub>50</sub>NO<sub>7</sub><sup>+</sup>

# *Vision for the future*

and is a form of data reuse











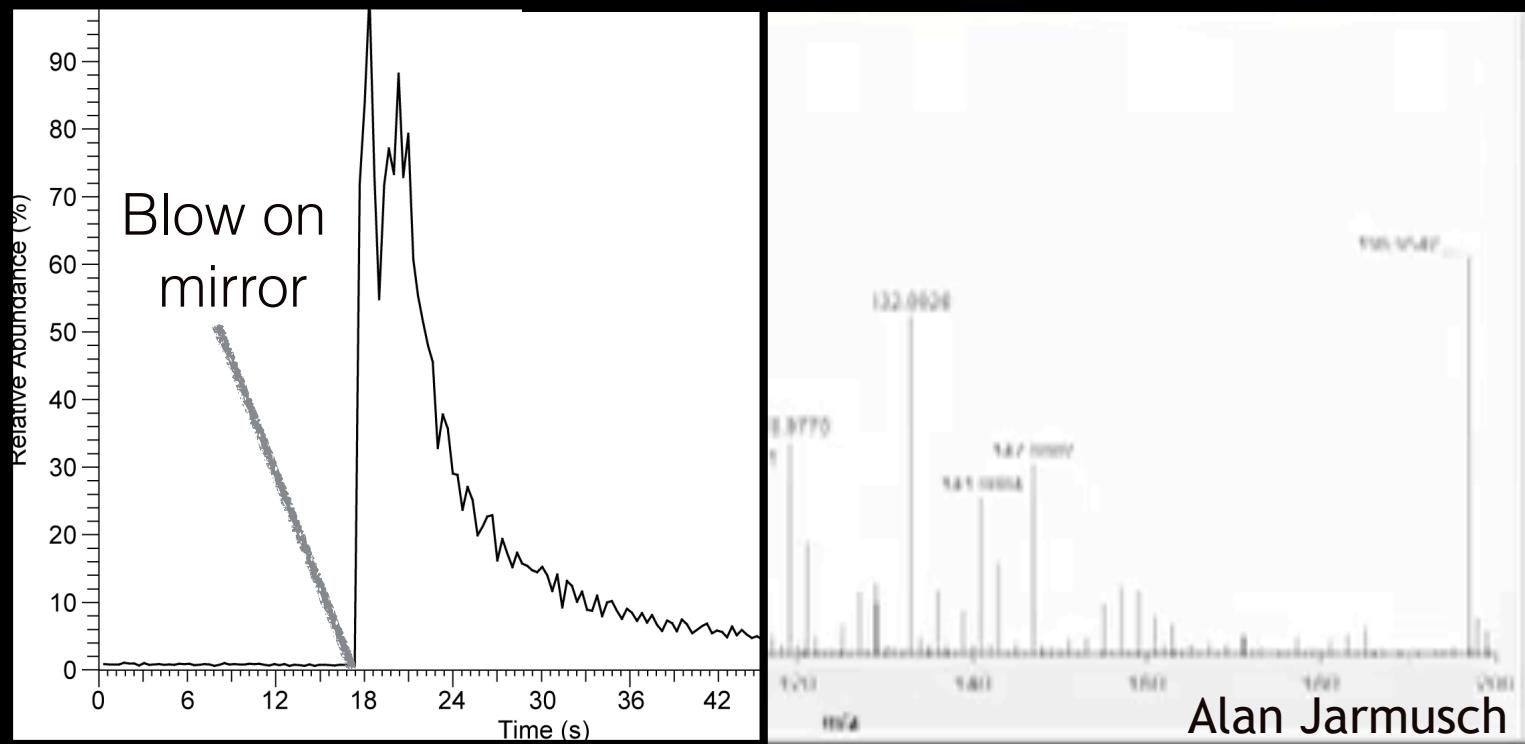


DANGER:  
IBD RISK!



**Go and see you doctor!!**

# Mirror Mass Spectrometry

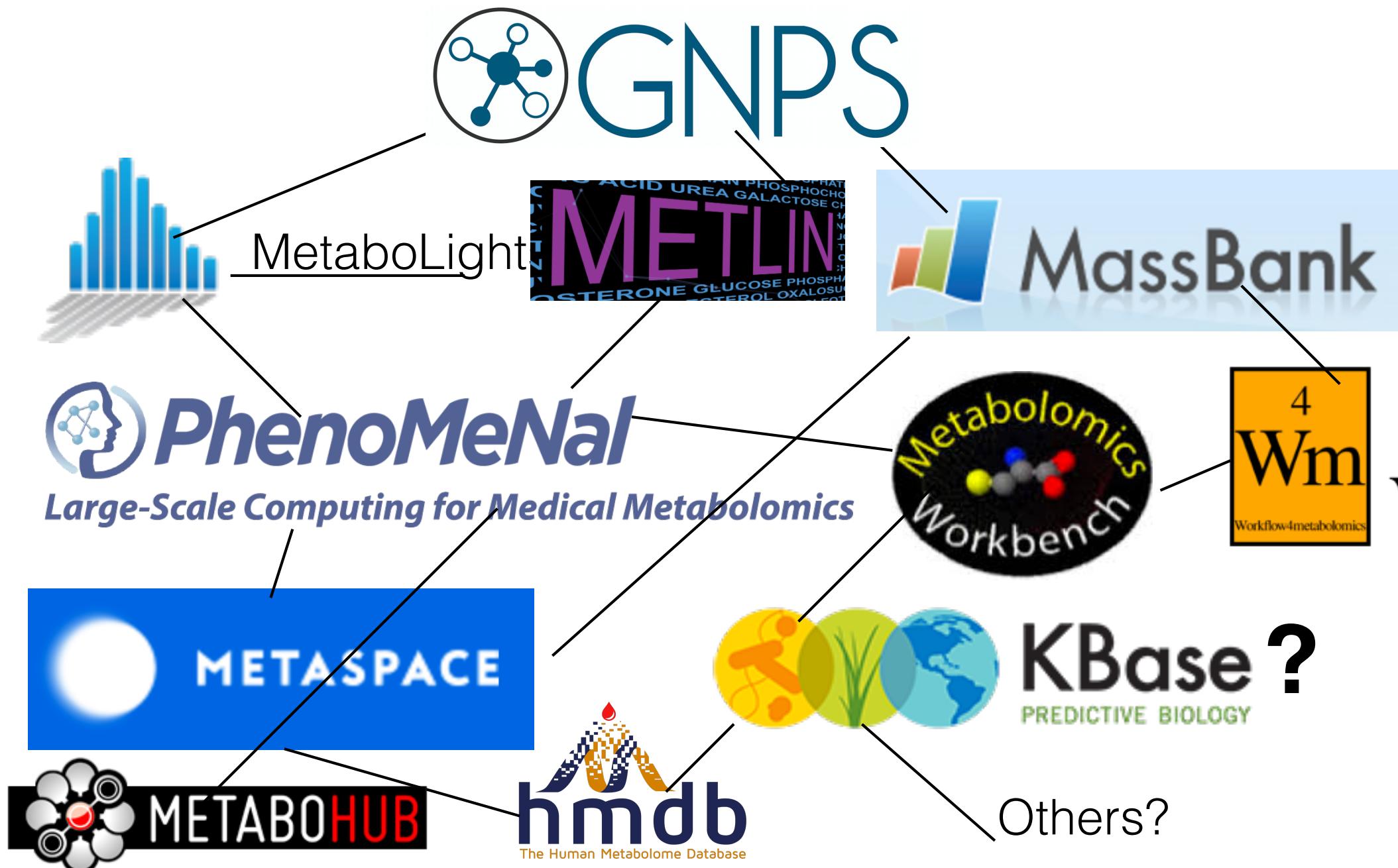


Alan Jarmusch

Data availability and algorithms are the limitation



Ultimately a network of interoperation need to be established to analyze MS data- as they all serve different functions. What might this look like?



Ultimately a network of interoperation need to be established to analyze MS data- as they all serve different functions.

The screenshot shows the homepage of **nature biotechnology**. The header features a green banner with the journal name and a background image of a coral reef and a circuit board. Below the banner, there's a blue navigation bar with a menu icon and a link to the **ARTICLE PREVIEW**. The main content area displays the title **SPLASH, a hashed identifier for mass spectra** by a large team of authors, including Gert Wohlgemuth, Sajjan S Mehta, Ramon F Mejia, Steffen Neumann, Diego Pedrosa, Tomáš Pluskal, Emma L Schymanski, Egon L Willighagen, Michael Wilson, David S Wishart, Masanori Arita, Pieter C Dorrestein, Nuno Bandeira, Mingxun Wang, Tobias Schulze, Reza M Salek, Christoph Steinbeck, Venkata Chandrasekhar Nainala, Robert Mistrik, Takaaki Nishioka & Oliver Fiehn. The footer includes the **Human Metabolome Database**.

**Large-S**

4  
Wm  
Workflow4metabolomics

?

The Human Metabolome Database

1) Raw data has to be shared-this is key.

## Take home messages

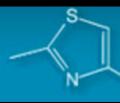
- 2) If the same data is used, and same tools with the same time stamp, one should get the same data. In other words cloning of the exact analysis of the method with the same parameters is possible. Metadata collection is critical and will need to be standardized
- 3) Multiple resources should communicate-data standards enables such cross communication.

nature  
chemical biology

## Large scale dereplication

NATURE CHEMICAL BIOLOGY

CURRENT ISSUE // ARCHIVE // NEWS & MULTIMEDIA // AUTHORS // ABOUT



## Computational development

Searching molecular structure database spectra using CSIL:Fin

Kai Dührkop<sup>a</sup>, Huibin Shen<sup>b</sup>, Marvin

Trends in Pharmacological Sciences

## Precision medicine

2017

Molecular Networking As a Drug Discovery, Drug Metabolism, and Proceedings of the National Academy of Sciences of the United States of America Precision Medicine Strategy

// AUTHORS // ABOUT COLLECTED ARTICLES

Dereplication  
database search

Hosein Mohimani, Alex  
Akihiro Ninomiya, Kent

## Forensics

Lifestyle chemistries from phones for individual

mohimani, E7645–E7654, doi: 10.1073/pnas.1610010114

## Microbial interactions

NATURAL  
PRODUCTS

Revisiting Previously Investigated Plants: A Molecular Networking-Based Study of *Geissospermum laeve*

Alexander E. Fox Ramos<sup>†</sup>, Charlotte Alcover<sup>†</sup>, Laurent Evanno<sup>†</sup>, Alexandre Maciuk<sup>†</sup>, Marc Litaudon<sup>‡</sup>, Christophe Duplais<sup>§</sup>, Guillermo Julian<sup>†</sup>, Elisabeth Mouray<sup>†</sup>, Philippe Gré Champy<sup>†</sup>, and Mehdi A. Benidir<sup>†</sup>

Proceedings of the National Academy of Sciences of the United States of America

## MS based Genome mining

Comparative genomics uncovers the prolific and distinctive metabolic potential of the cyanobacterial genus *Moorea*

Tiago Leao<sup>a</sup>, Guilherme Castelão<sup>b</sup>, Anton Korobeynikov<sup>c,d</sup>, Emily A. Monroe<sup>e</sup>, Sheila Podell<sup>a</sup>,

Evguenia Glukhov<sup>a</sup>, Eric F. Allen<sup>a</sup>, William H. Gerwick<sup>a,f</sup>, and Lena Gerwick<sup>a,1</sup>

## Reuse of old data

Expanding the Chemical Repertoire of the Endophyte *Streptomyces albospinus* RLe7 Reveals Amphotericin B as an Inducer of a Fungal Phenotype

Mauricio Caraballo-Rodríguez<sup>†,‡</sup>, Pieter C. de Groot<sup>†</sup>,

ECOLOGY  
ECOLOGICAL SOCIE

## Forest ecology

Sources of variation in foliar secondary chemistry in a tropical forest tree community

Brian E. Sedio<sup>✉</sup>, Juan C. Rojas Echeverri, Cristopher A. Boya P., S. Joseph Wright

Ocean DOM

## Exposomics

Sci Rep. 2017; 7:  
published online 2017 Mar 7. doi: 10.1038/srep44038

In vivo microsampling to capture the elusive exposome

Vincent Bessonneau<sup>1</sup>, Jennifer Ings<sup>2</sup>, Mark McMaster<sup>2</sup>, Richard Smith<sup>3</sup>, Leslie Bragg<sup>4</sup>, Mark Servos<sup>4</sup>, and Janusz Pawliszyn<sup>a,1</sup>

# THANKS!

## The community to crowdsource molecular networking

Mingxun Wang<sup>1,3</sup>, Jeremy Carver<sup>1,3</sup>, Vanessa V. Phelan<sup>1,5</sup>, Laura M. Sanchez<sup>1,5</sup>, Neha Garg<sup>1,5</sup>, Yao Peng<sup>1,4</sup>, Jeramie Watrous<sup>5</sup>, Clifford Kopano<sup>4</sup>, Don Duy Nguyen<sup>4</sup>, Tal Luzzatto-Knaan<sup>5</sup>, Carla Porto<sup>5</sup>, Amina Bouslimani<sup>5</sup>, Alexey V. Melnik<sup>5</sup>, Michael J. Meehan<sup>5</sup>, Wei-Ting Liu<sup>6</sup>, Paul D. Boudreau<sup>8</sup>, Eduardo Esquenazi<sup>10</sup>, Mario Sandoval-Calderón<sup>11</sup>, Roland D. Kersten<sup>12</sup>, Robert A. Quinn<sup>13</sup>, Katherine Duncan<sup>8,39</sup>, Richard Hsu<sup>4,6</sup>, Dimitrios J. Floros<sup>4,5</sup>, Ronnie G. Gavilan<sup>14</sup>, Karin Kleigrewe<sup>8</sup>, Norberto Peoporine Lopes<sup>16</sup>, Trent Northen<sup>17</sup>, Rachel J. Dutton<sup>46</sup>, Delphine Parrot<sup>19</sup>, Erin E. Carlson<sup>30</sup>, Bertrand Aigle<sup>31</sup>, Lars Jelsbak<sup>29</sup>, Christian Sohlenkamp<sup>11</sup>, Charlotte F. Michelsen<sup>29</sup>, Pavel Pevzner<sup>3</sup>, Anna Edlund<sup>33, 34</sup>, Jeffrey McLean<sup>28</sup>, Joern Piel<sup>26</sup>, Brian T. Murphy<sup>27</sup>, Lena Gerwick<sup>8</sup>, Chih-Chuang Liaw<sup>22</sup>, Yu-Liang Yang<sup>23</sup>, Hans-Ulrich Humpf<sup>40</sup>, Maria Maansson<sup>24</sup>, Robert A Keyzers<sup>25</sup>, Amy C Sims<sup>36</sup>, Andrew R. Johnson<sup>37</sup>, Ashley M. Sidebottom<sup>37</sup>, Brian E. Sedio<sup>38</sup>, Andreas Klitgaard<sup>29</sup>, Charles B. Larson<sup>8,42,49</sup>, Christopher A. Boya P.<sup>14</sup>, Daniel Torres-Mendoza<sup>14</sup>, David J. Gonzalez<sup>5</sup>, Denise Brentan Silva<sup>16</sup>, Egle Pociute<sup>10</sup>, Ellis O'Neill<sup>8</sup>, Enora Briand<sup>8, 19</sup>, Eric J. N. Helfrich<sup>26</sup>, Eve A. Granatosky<sup>32</sup>, Evgenia Glukhov<sup>8</sup>, Florian Ryffel<sup>26</sup>, Hailey Houson<sup>10</sup>, Hosein Mohimani<sup>3</sup>, Jenan Kharbush<sup>8</sup>, Yi Zeng<sup>4</sup>, Julia A. Vorholt<sup>26</sup>, Kenji L. Kurita<sup>21</sup>, Pep Charusanti<sup>15</sup>, Kerry L. McPhail<sup>44</sup>, Kristian Fog Nielsen<sup>29</sup>, Lisa Vuong<sup>10</sup>, Maryam Elfeki<sup>27</sup>, Matthew F. Traxler<sup>18</sup>, Max Crüsemann<sup>8</sup>, Niclas Engene<sup>47</sup>, Nobuhiro Koyama<sup>5</sup>, Oliver B. Vining<sup>44</sup>, Ralph Baric<sup>36</sup>, Ricardo Roberto Silva<sup>16</sup>, Samantha J Mascuch<sup>8</sup>, Sophie Tomasi<sup>19</sup>, Stefan Jenkins<sup>17</sup>, Venkat Macherla<sup>10</sup>, Thomas Hoffmann<sup>51</sup>, Vinayak Agarwal<sup>52</sup>, Philip G. Williams<sup>54</sup>, Jingqui Dai<sup>54</sup>, Ram Neupane<sup>54</sup>, Joshua Gurr<sup>54</sup>, Andrés Mauricio Caraballo Rodríguez<sup>53</sup>, Wenyuan Shi<sup>50</sup>, Rob Knight<sup>41</sup>, Paul R. Jensen<sup>8</sup>, Bernhard Ø. Palsson<sup>15</sup>, Kit Pogliano<sup>7</sup>, Roger G. Linington<sup>21</sup>, Marcelino Gutiérrez<sup>14</sup>, William H. Gerwick<sup>5,8</sup>, Bradley S. Moore<sup>5,8,52</sup>, Nuno Bandeira<sup>2,3,5</sup>

# A Big Thanks to Friends, Co-workers and Support That Allow Us To Do The Research

## Bandeira Lab

Ming Wang

## Alexandrov lab

Sergey Ryazanov

Ivan Protsyuk

## Knight lab

Antonio Peña,  
Donna Berg-Lyons

Gail Ackermann

Embriette Hyde

Gregory Humphrey

James Gaffney

Jeff Dereus

Dominguez lab

Pilau lab

## Moore Lab

## Gerwick Labs

## Jensen lab

## Fenical Lab

## Nizet Lab

## Pevzner lab

## Rohwer lab

## Conrad lab

## Lin lab

## Medema lab

## Boeker lab

## Mazmanian lab

## Dorrestein lab (current)

Neha Garg

Robby Quinn

Cliff Kapono

Fernando Vargas

Don Nguyen

Amina Bouslimani

Alexey Melnik

Kathleen Dorrestein

Mike Meehan

Carla Porto

Ricardo da Silva

Alexander Aksenov

Alan Jarmusch

Julia Gauglitz

And visitors

Thanks to reviewers  
of grants that have  
supported forward  
looking ideas

**Current Support:**  
Bruker, NIH, NSF, USDA  
EU 7<sup>th</sup> framework  
and horizon 2020  
Janssen, Bayer, Colgate  
Sloan and Moore  
foundations  
UCSD, Department  
of Justice, Office of  
Naval Research, Science  
without Borders,  
lululemon, CMI

Other collaborators and co-workers were highlighted throughout the talk





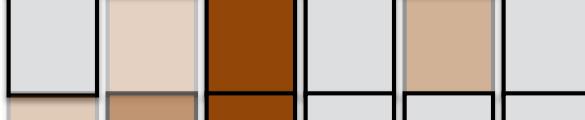
**Activity**  
(EC50 in  $\mu\text{g}/\text{ml}$ )

*m/z* A B C D E F

Fraction 1



Fraction 2



Fraction 3



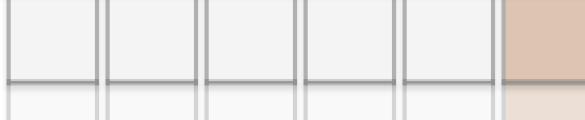
Fraction 4



Fraction 5

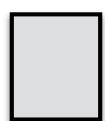
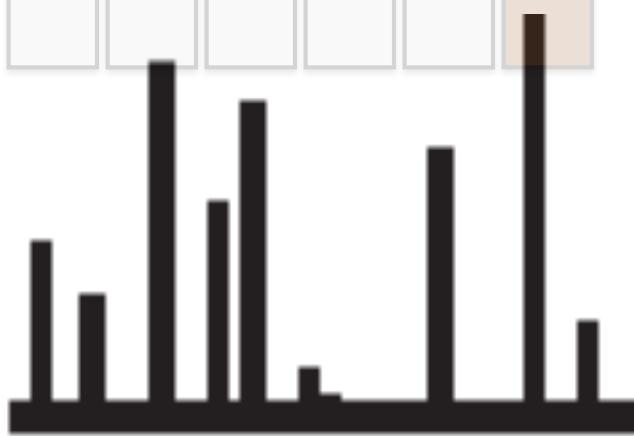


Fraction 6

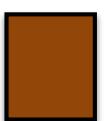


Fract....

Etc.

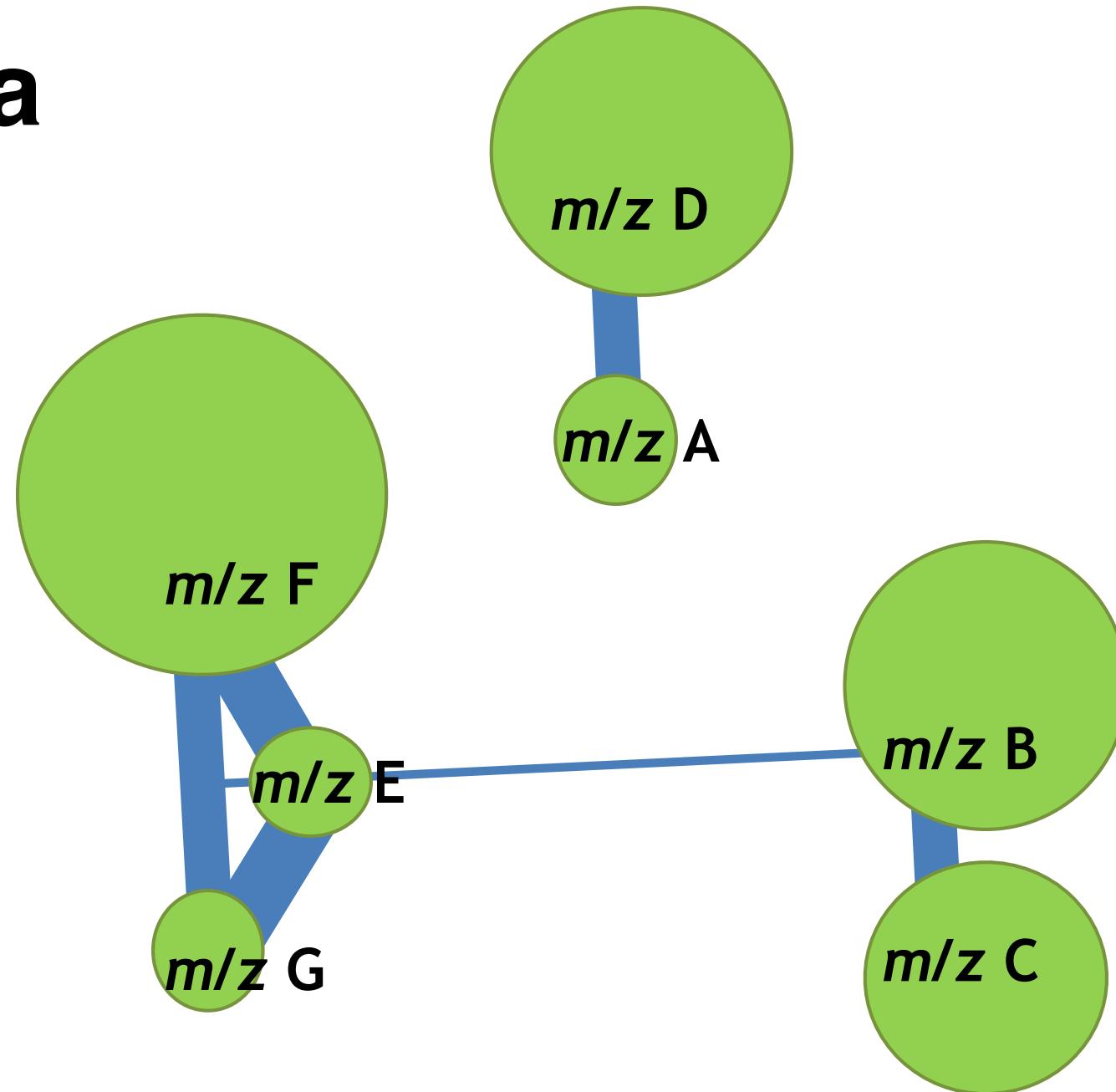


= not detected/no activity



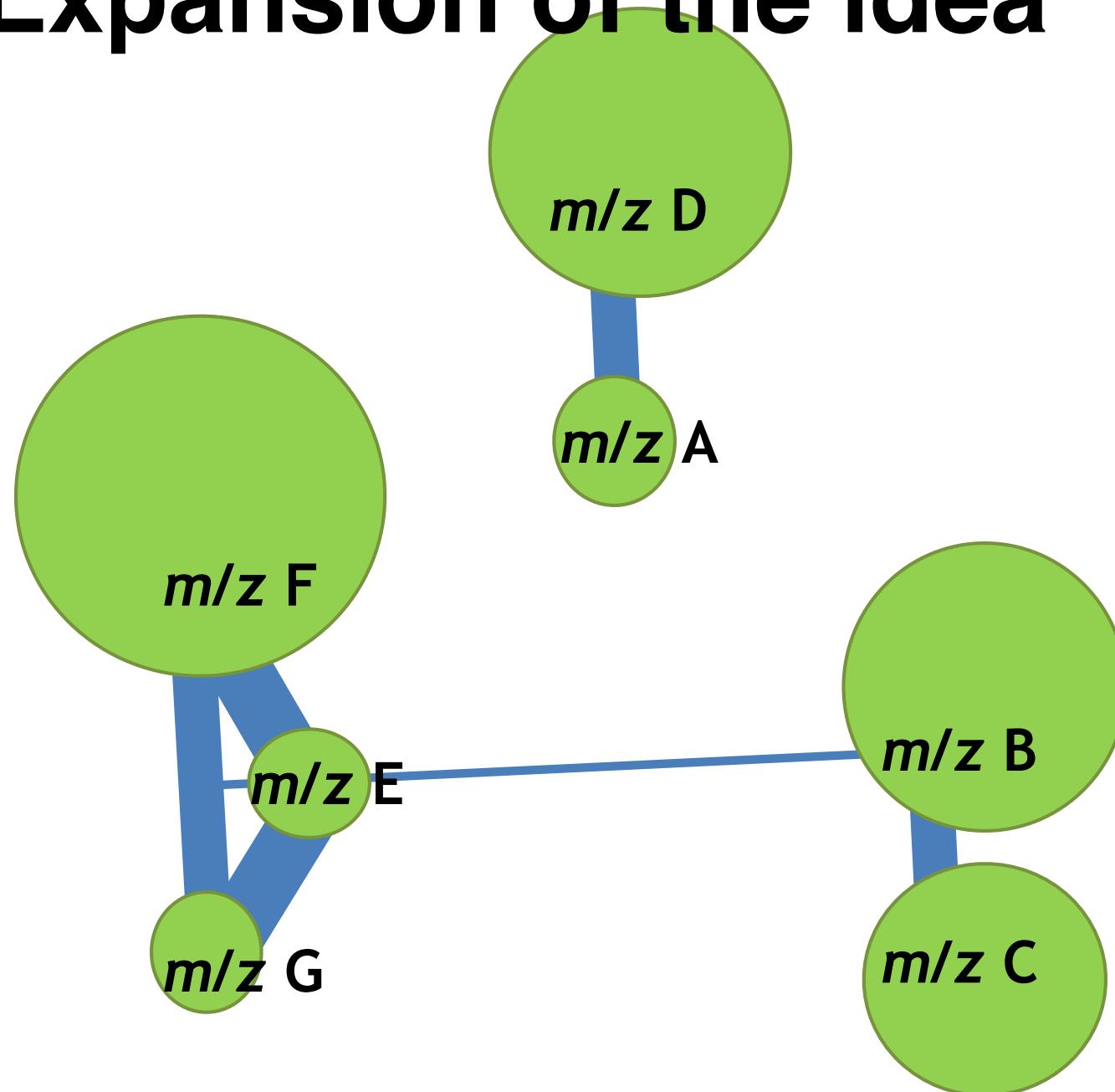
= most intense/active

# Idea



Problem - many ions observed in the three active fractions

# Expansion of the idea



## LEGEND

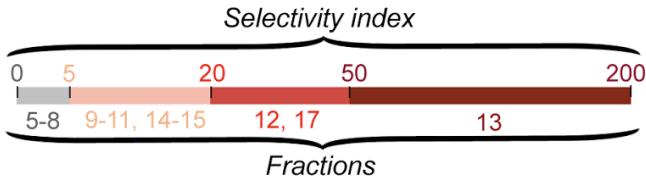
■ Compounds identified by comparison with standards (1-19) or isolated (20-23)

Bioactivity score:  Bioactivity score:

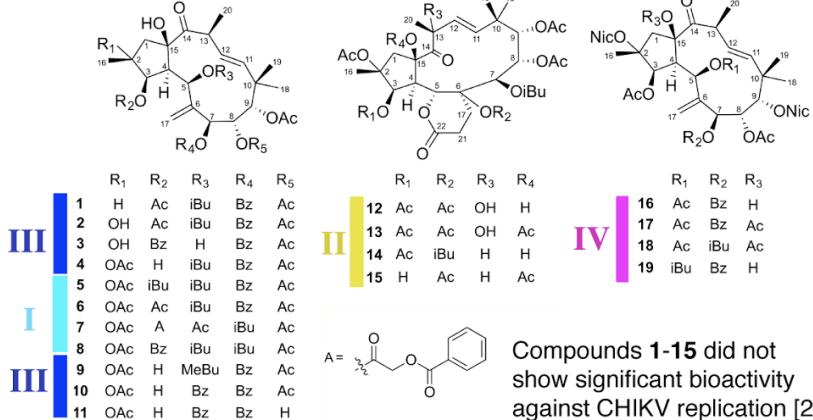
Node sizes:   $r^2 < 0.85$    $r^2 > 0.85$

p-value > 0.03  p-value < 0.03

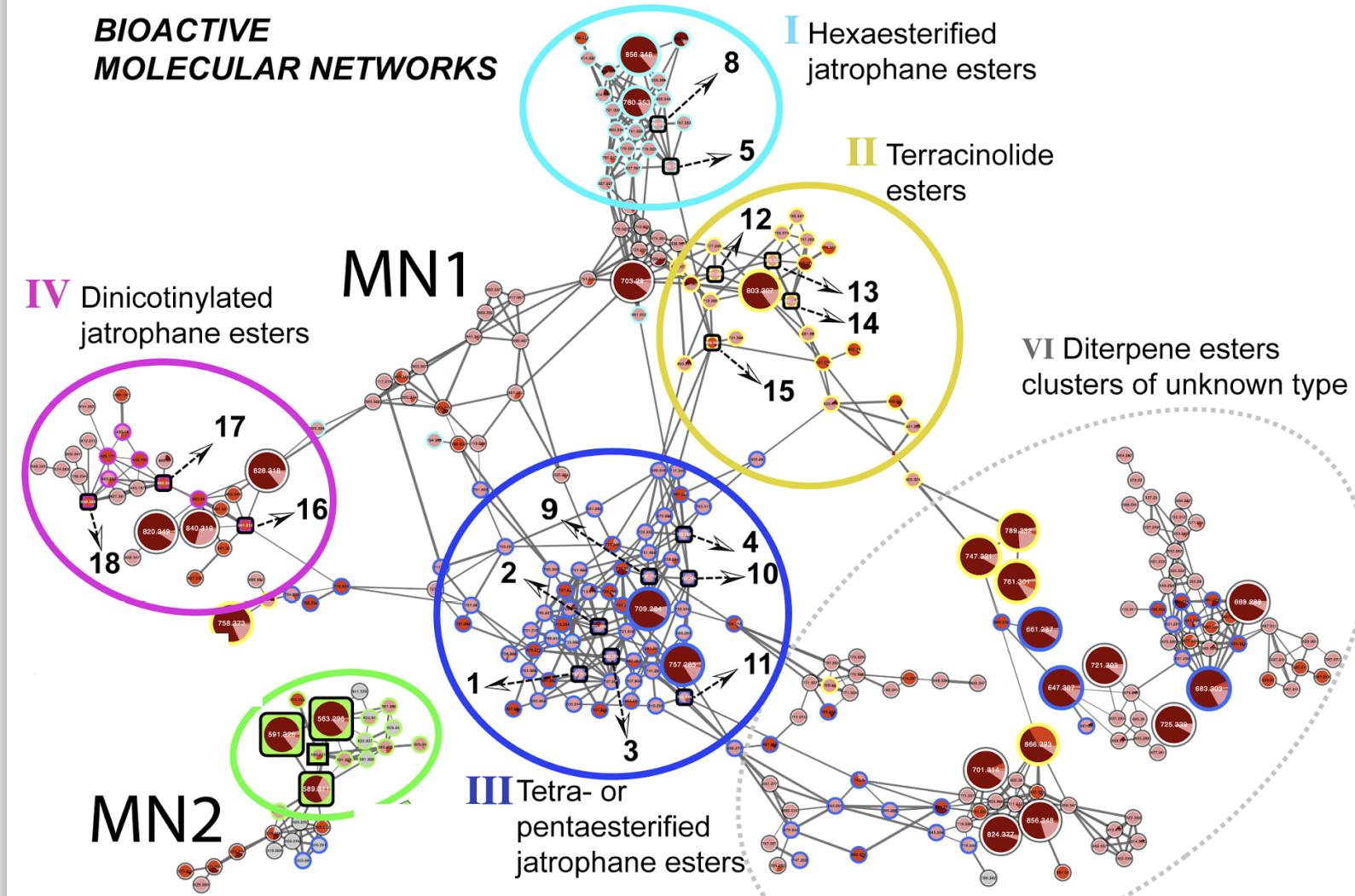
● Relative quantification of the feature accross fractions



## PREVIOUSLY ISOLATED COMPOUNDS (1-19)



## BIOACTIVE MOLECULAR NETWORKS



# Natural products from food

