

Improved small molecule structural annotation by exploiting biochemical and biosynthetic relationships

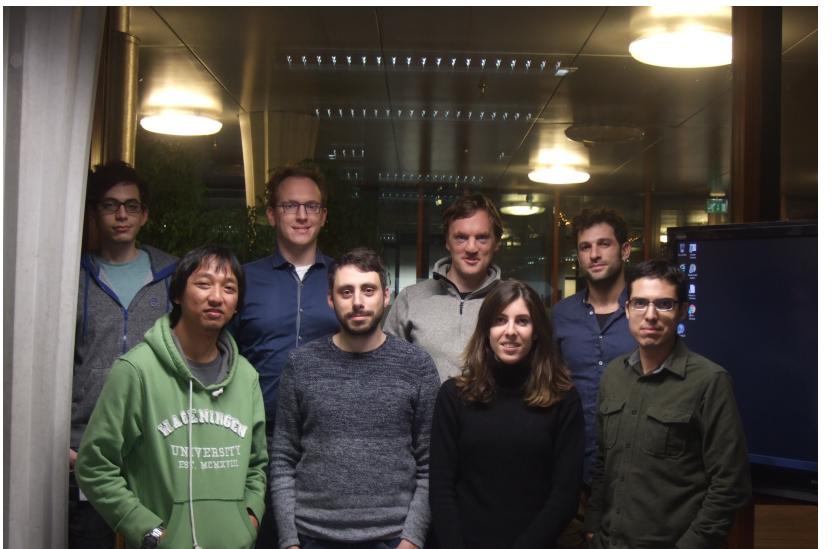
Justin J.J. van der Hooft et al.

Bioinformatics Group – Wageningen University, The Netherlands

Antwerpen, 30 November 2018



Team work! ☺



Medema lab - Wageningen UR, NL



NL eScience Center



Dorrestein lab – San Diego, USA



Glasgow Polyomics – University of Glasgow, UK



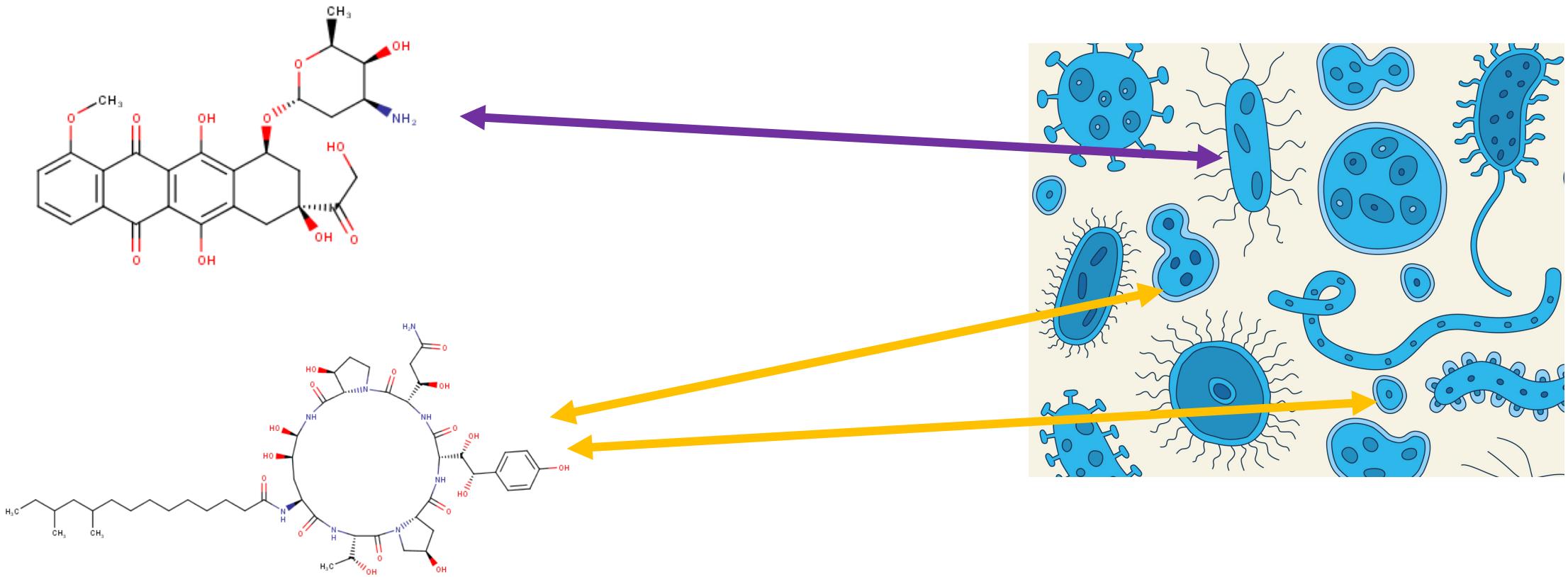
€€ Funding €€ **NWO eScience center**



Glasgow Polyomics
www.glasgow.ac.uk/polyomics

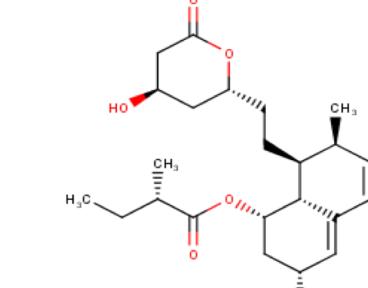
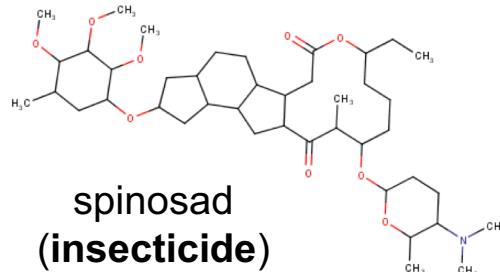


iOMEGA project: integrated omics for metabolomics and genomics annotation

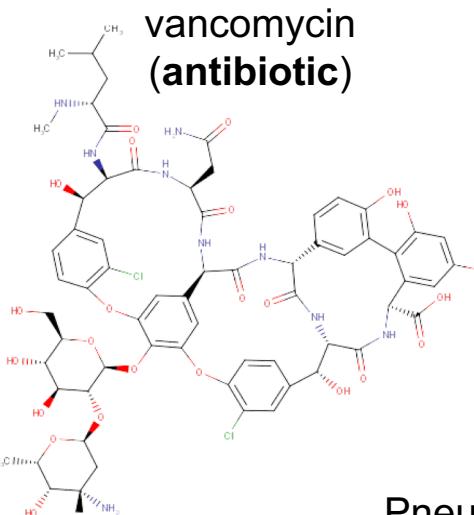
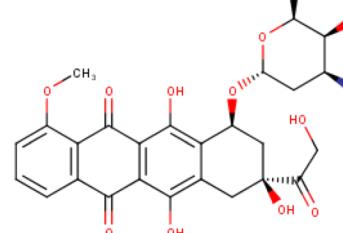


The challenge....

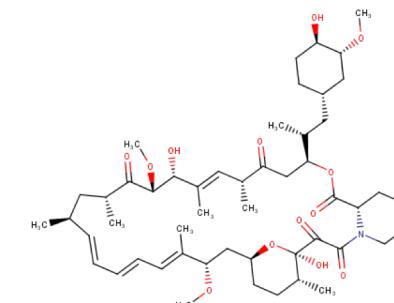
Bacteria, fungi, and plants produce a large & diverse arsenal of high-value molecules:



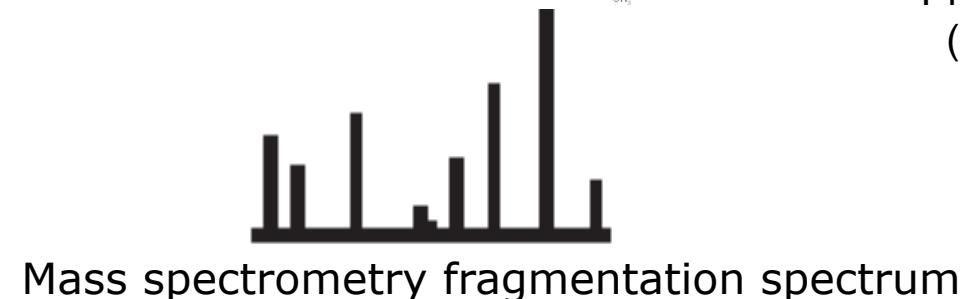
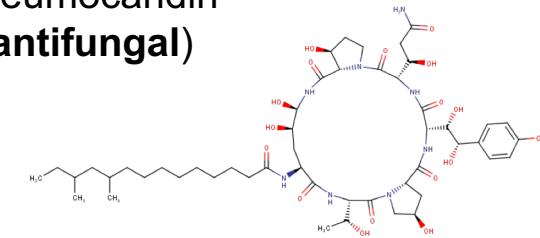
doxorubicin
(chemotherapeutic
agent)



rapamycin
(immunosuppressant)



Pneumocandin
(antifungal)



....is large-scale coupling of spectral data to molecular structures
of known & especially **novel** natural products molecules

Motivation



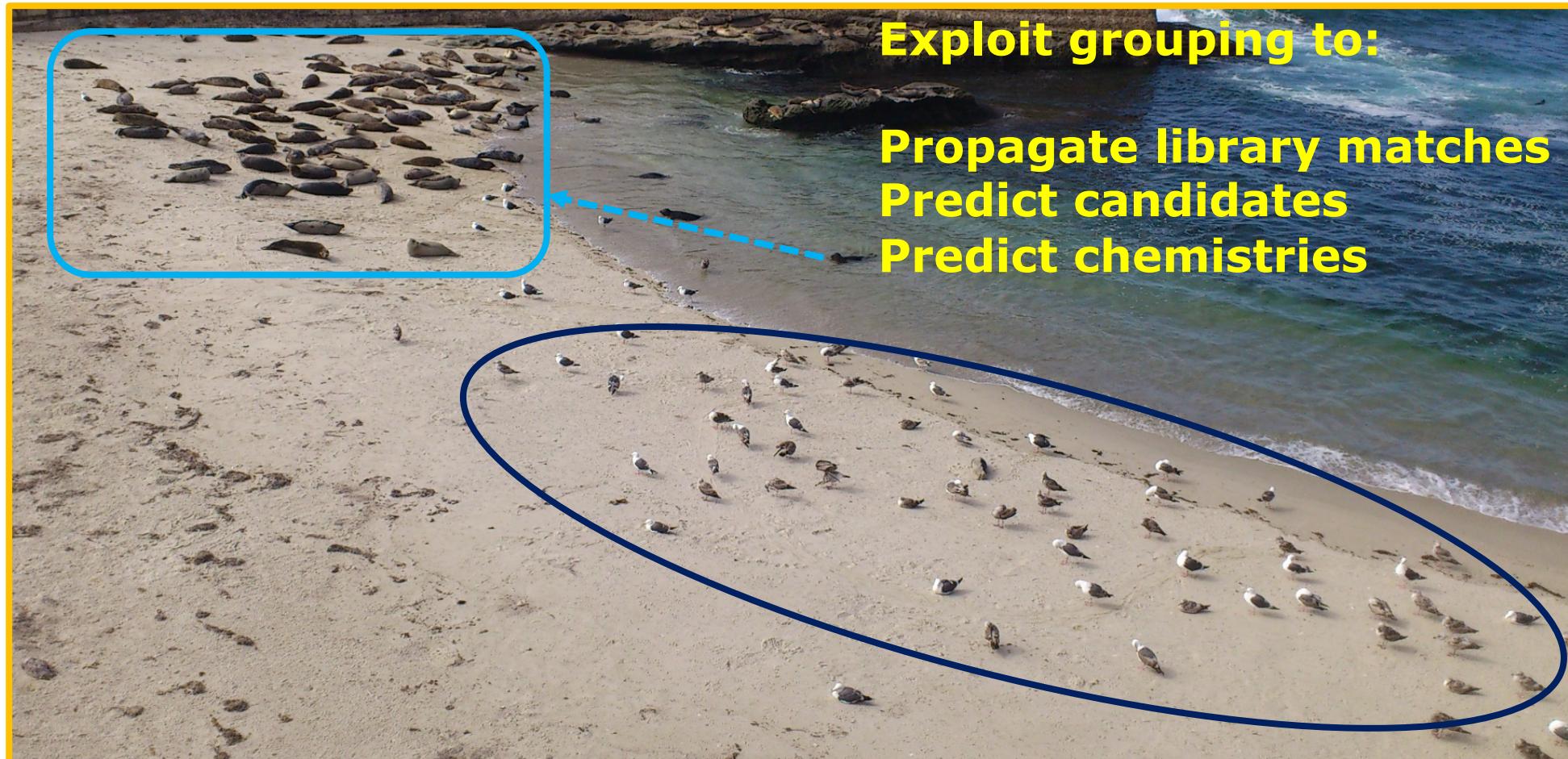
MS/MS information is key to study large sample sets with diverse chemistries

Integrated workflow exploits complementary tools to enhance interpretation

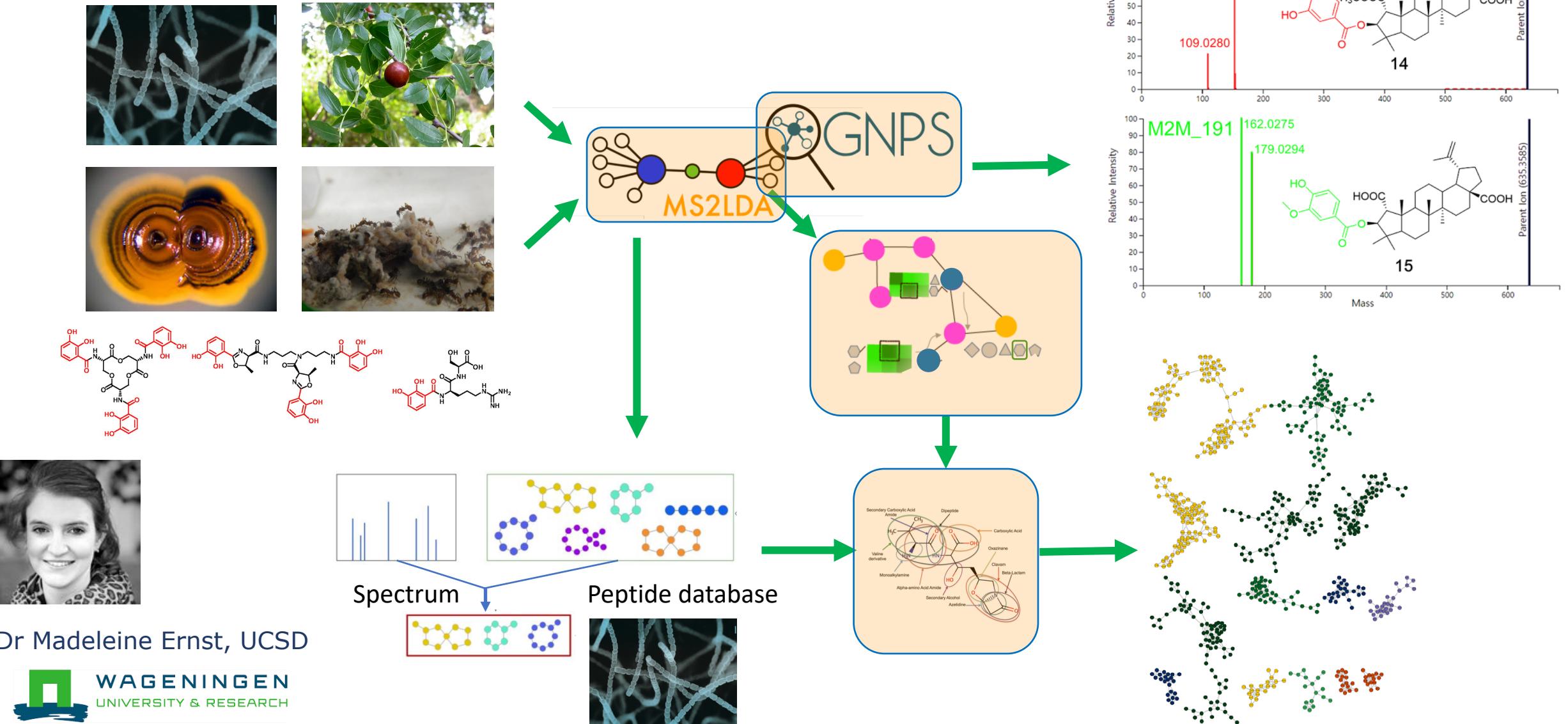
Motivation: Improved annotation power by pattern mining

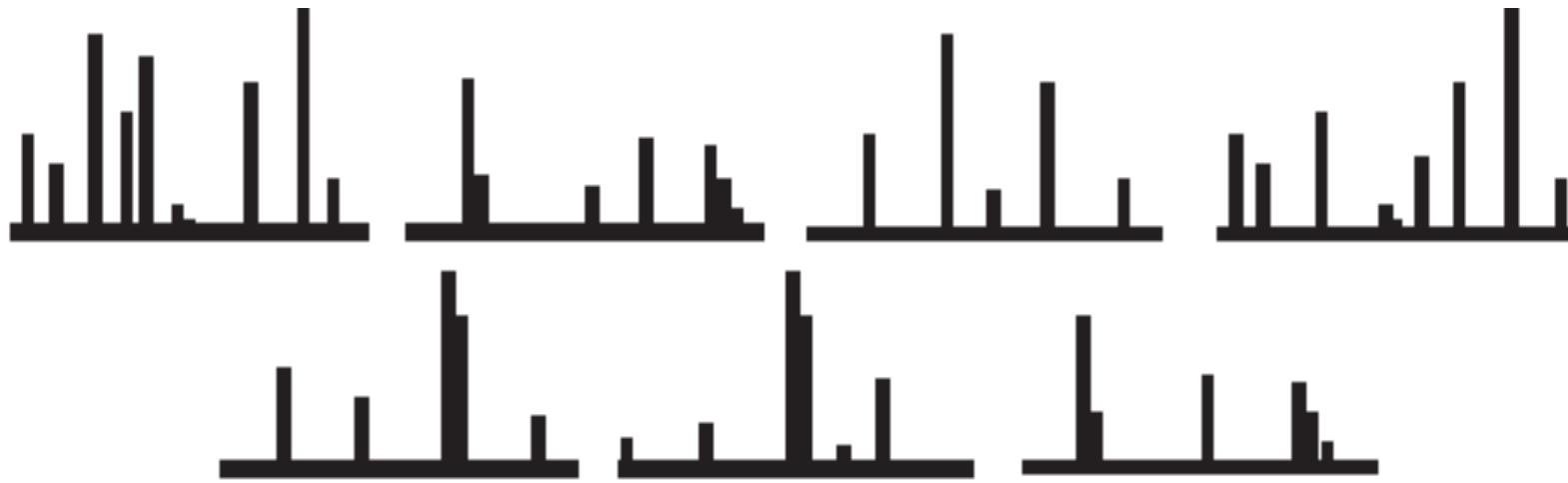
Molecular Families based on spectral similarity (Molecular Networking in GNPS)

Substructures by extracting “building blocks of metabolomics” (MS2LDA)



Integrated workflow





Molecular Networking

Very similar MS/MS spectra are grouped to:

- link spectra across different samples
- reduce redundancy in data set ("consensus spectrum")

Wang, M., et al., "Sharing and community curation of mass spectrometry data with Global Natural Products Social Molecular Networking"

Nat Biotech (2016)

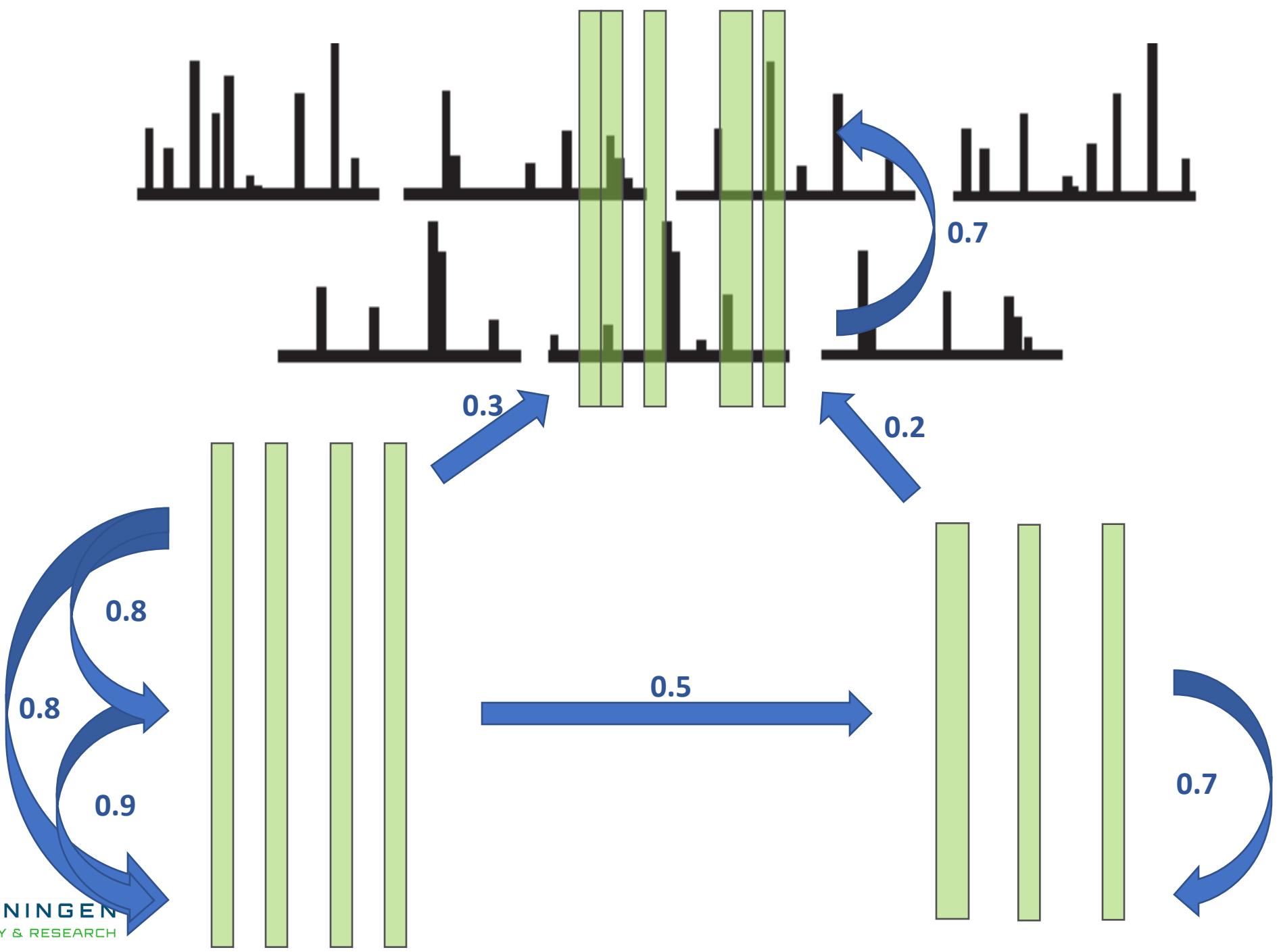
Watrous, JD et al. "Mass spectral molecular networking of living microbial colonies" PNAS (2012)

Dr Ricardo R. da Silva

Dr Madeleine Ernst

Dr Mingxun Wang

Dr Louis-Felix Nothias



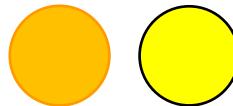
Spectral library matches from GNPS

Libraries from diverse sources

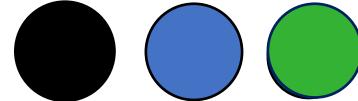
Seed node annotations for molecular families

Library MS/MS Spectra

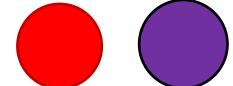
Diterpenoids



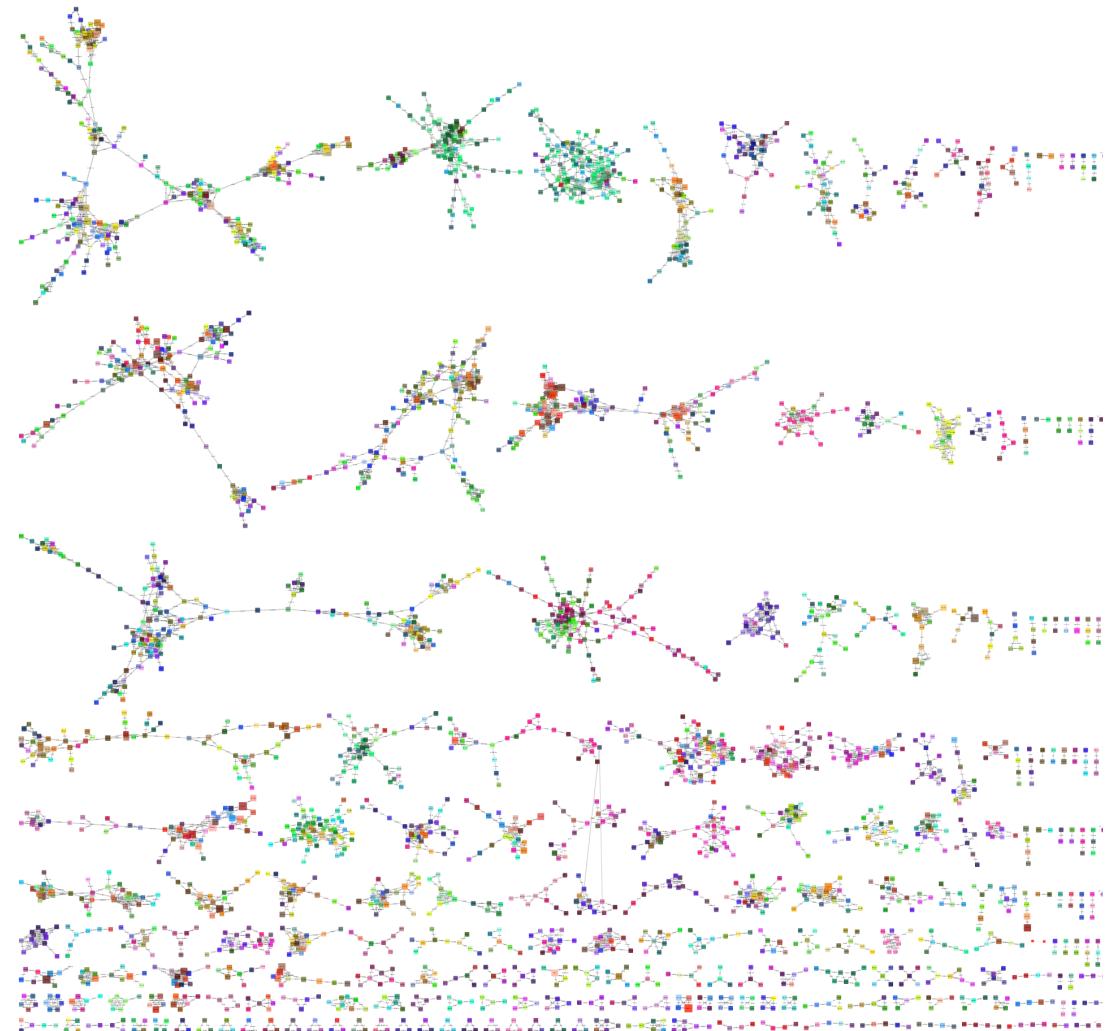
Flavonoids



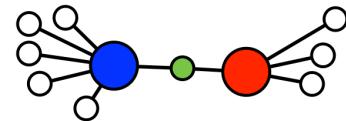
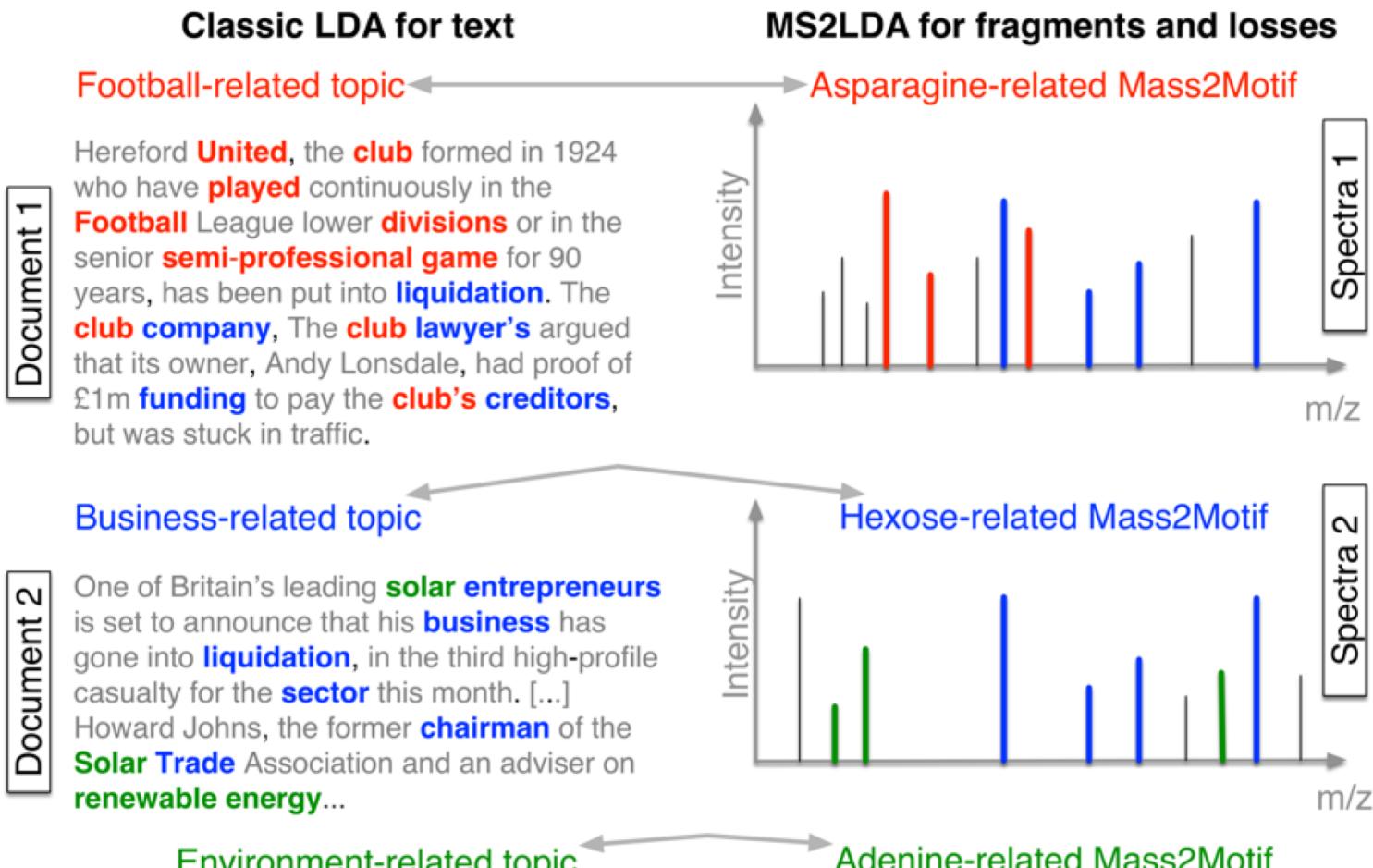
Pharmaceuticals



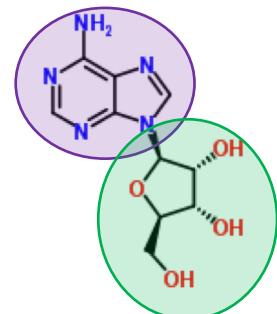
MS contaminants



Topic modelling: from text to molecules



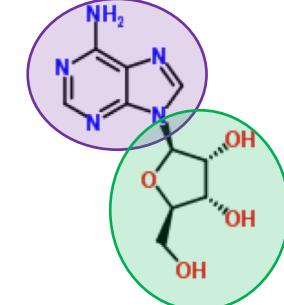
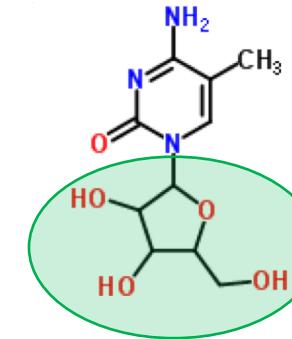
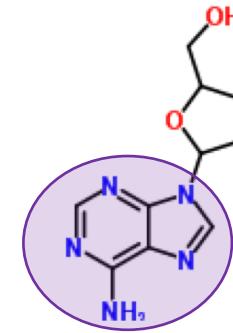
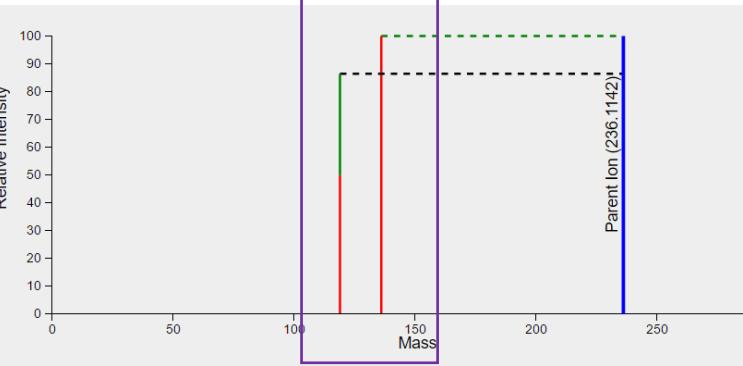
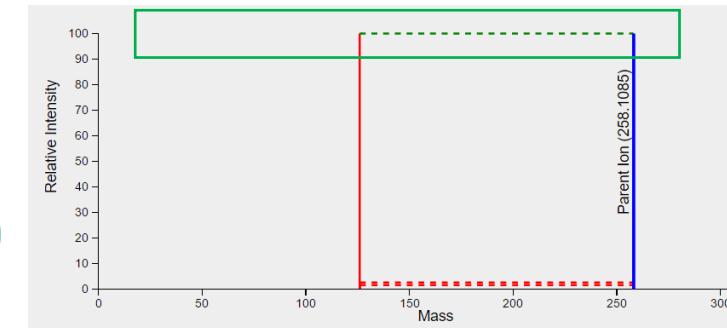
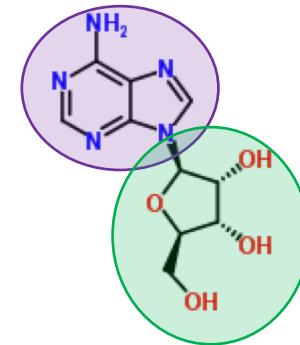
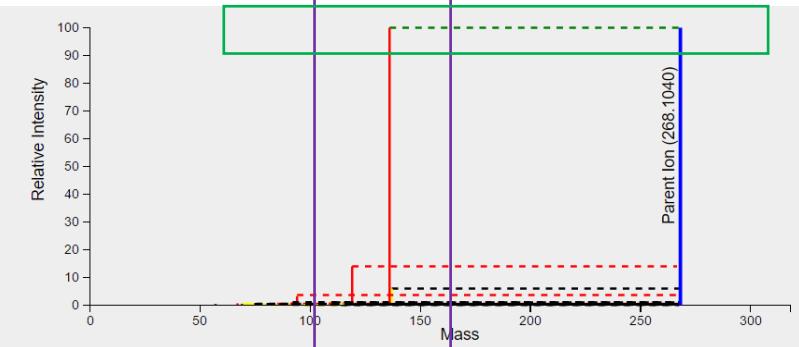
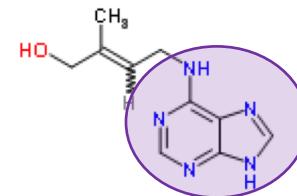
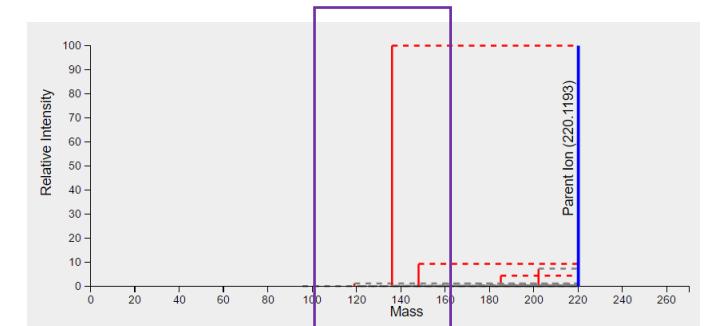
MS2LDA
Unsupervised Substructure Discovery



Documents \leftrightarrow molecules
Words \leftrightarrow fragments/neutral losses

Van der Hooft et al., PNAS, 2016

Validation: MS2LDA with standards

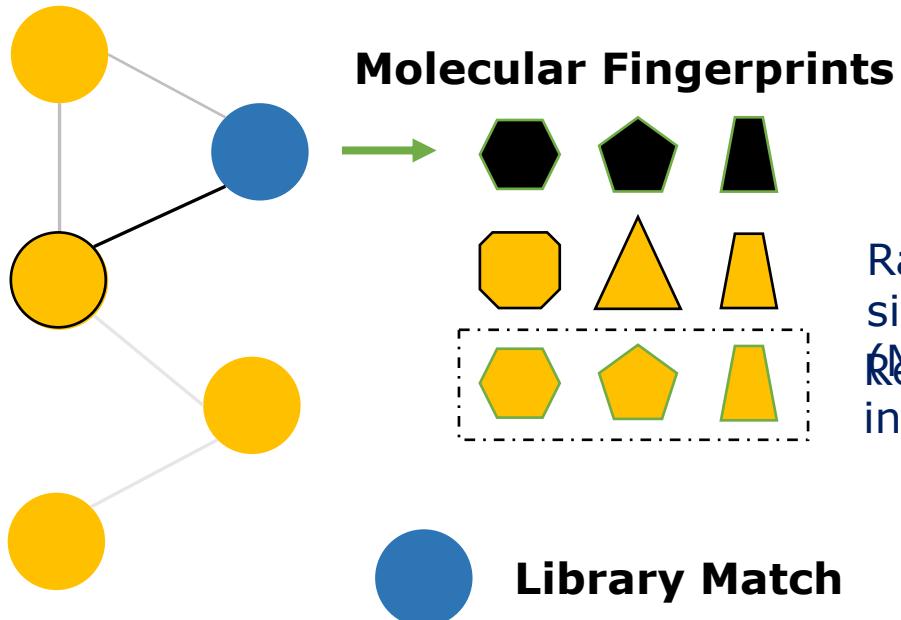


Network annotation propagation

More coherent candidate structures within family

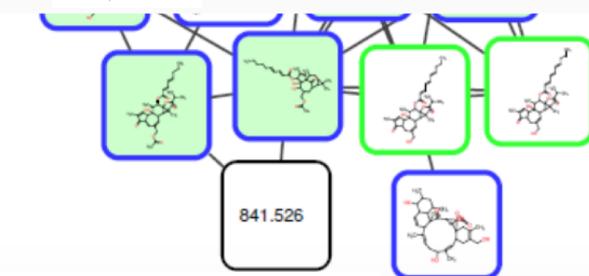
Rerank list based on neighboring annotations

Exploiting the network topology



Ranked lists of in silico annotators
~~(MetFrag)~~! lists of in silico annotators

| Scans | Assigned | Candidate score | Molecule | Formula | Mass | ΔMass (ppm) | Name |
|-------|----------|-----------------|----------|-------------|-----------|-------------|------------------------------------|
| 1 | No | 2.23513 | | C13H16N2O5S | 312.07799 | 0.86666 | acetaminophen mercapturate (83967) |
| 1 | No | 2.23513 | | C13H16N2O5S | 312.07799 | 0.86666 | (56923619) |
| 1 | No | 3.20109 | | C13H16N2O5S | 312.07799 | 0.86666 | AKOS019768418 (81197200) |
| 1 | No | 3.22166 | | C13H16N2O5S | 312.07799 | 0.86666 | CHEMBL1642089 (53321314) |
| 1 | No | 3.22166 | | C13H16N2O5S | 312.07799 | 0.86666 | CHEMBL1642098 (53318691) |
| 1 | No | 3.34424 | | C13H16N2O5S | 312.07799 | 0.86666 | AGN-PC-0G4R64 (64889340) |
| 1 | No | 3.42103 | | C13H16N2O5S | 312.07799 | 0.86666 | AGN-PC-0DZQ65 (61409875) |
| 1 | No | 3.42875 | | C13H16N2O5S | 312.07799 | 0.86666 | AGN-PC-0E3QTW (61600272) |



ClassyFire – SMILES and substituents

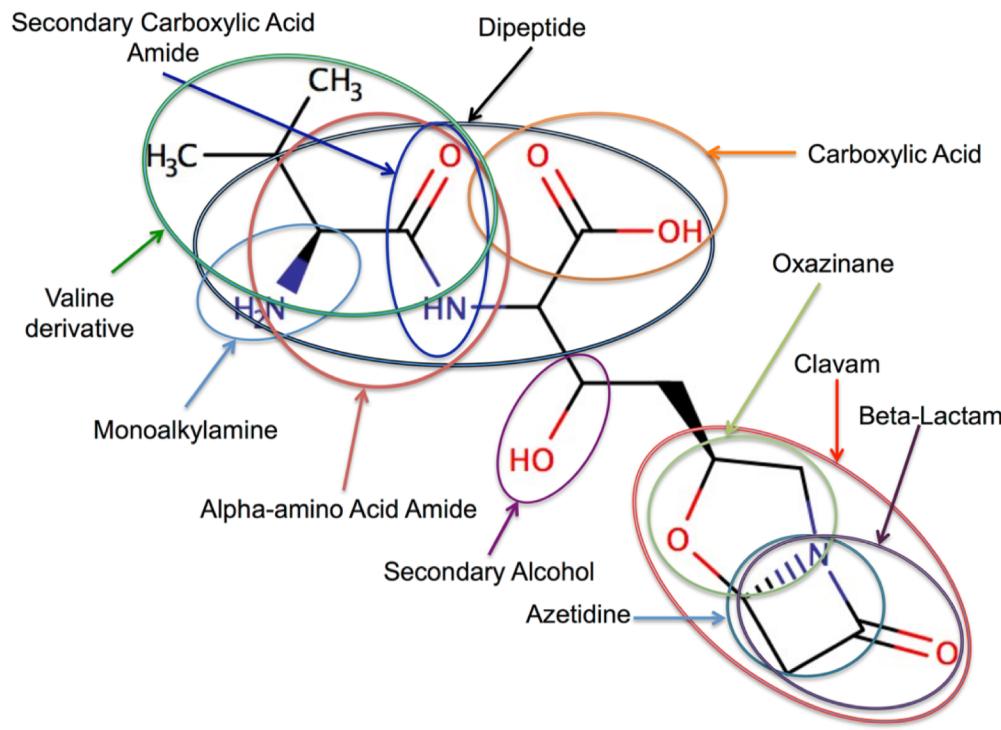


Figure part of Fig. 1 from Djoumbou Feunang et al., J. Cheminform, 2016

INPUT:

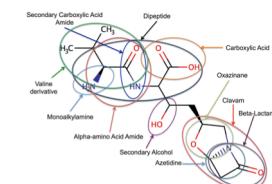
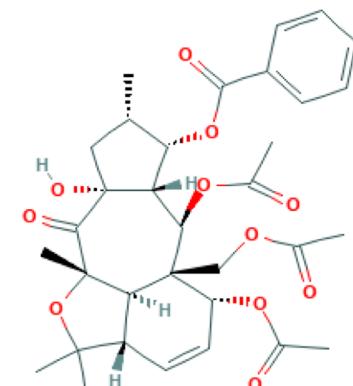
CC1CC2(C(C1OC(=O)C3=CC=CC=C3)C(C4(C(C=CC5C4C(C2=O)(OC5(C)C)C)OC(=O)C)COC(=O)C)OC(=O)C)O

OUTPUT:

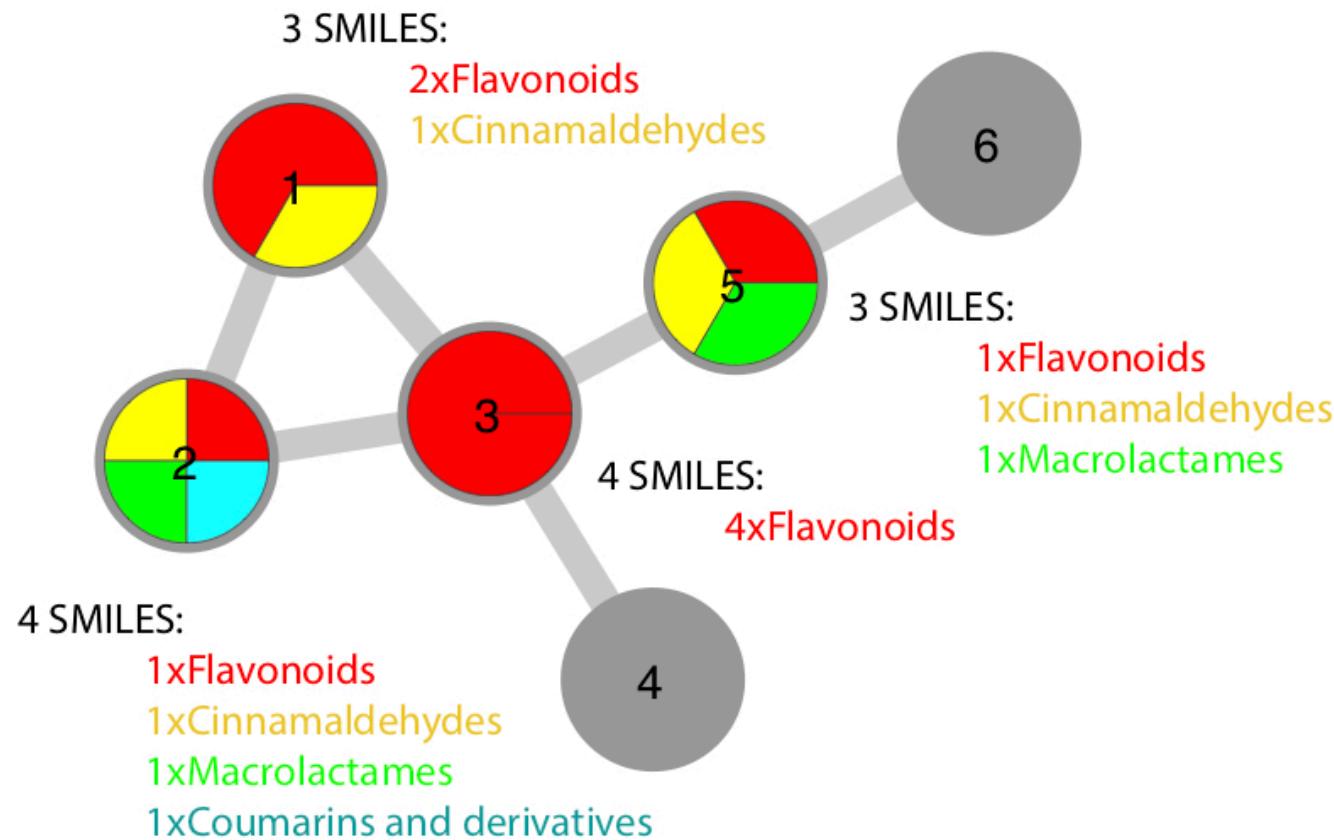
ClassyFire

Kingdom: Organic compounds
Superclass: Lipids and lipid-like molecules
Class: Prenol lipids
Subclass: Diterpenoids

...



ClassyFire – Chemical predictions for MFs



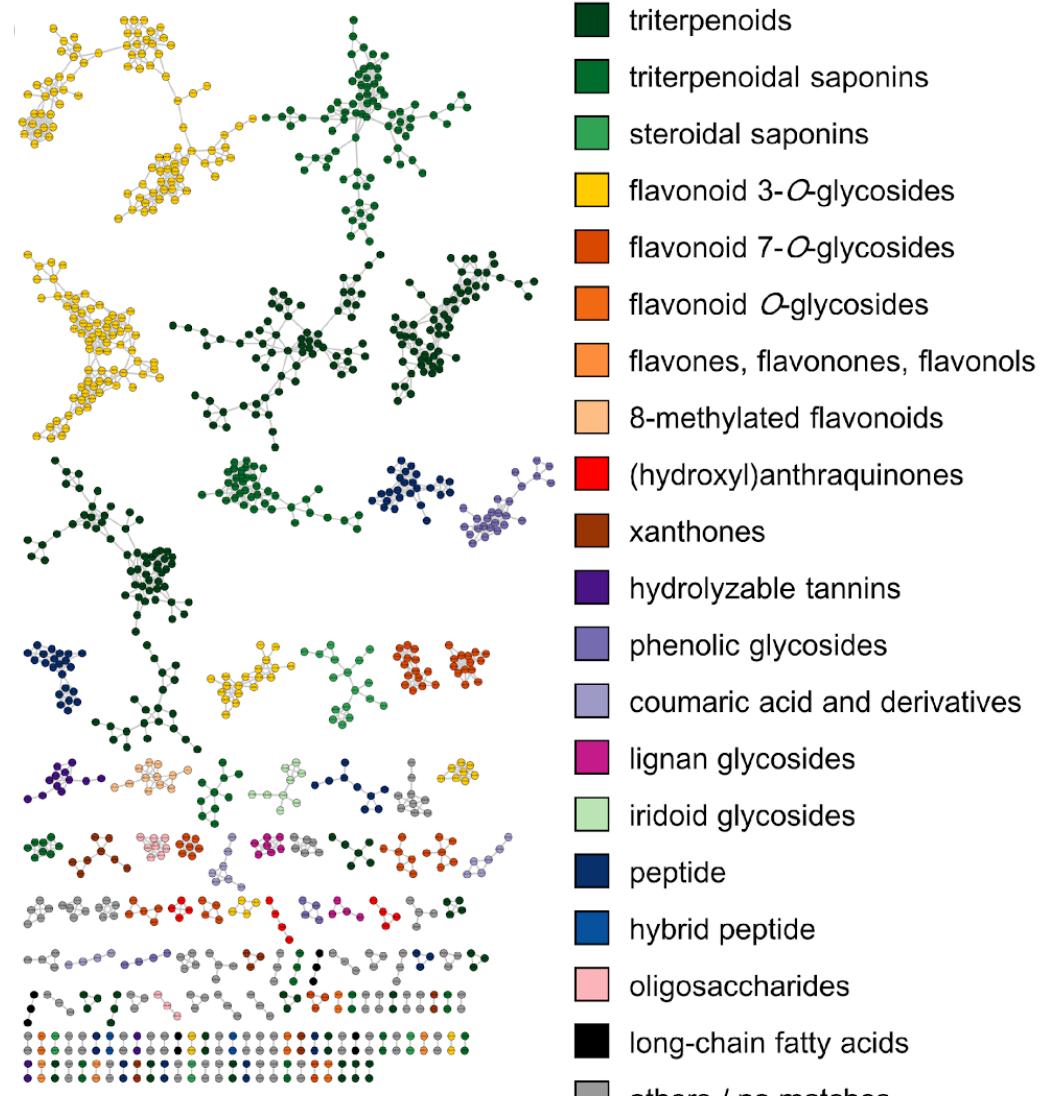
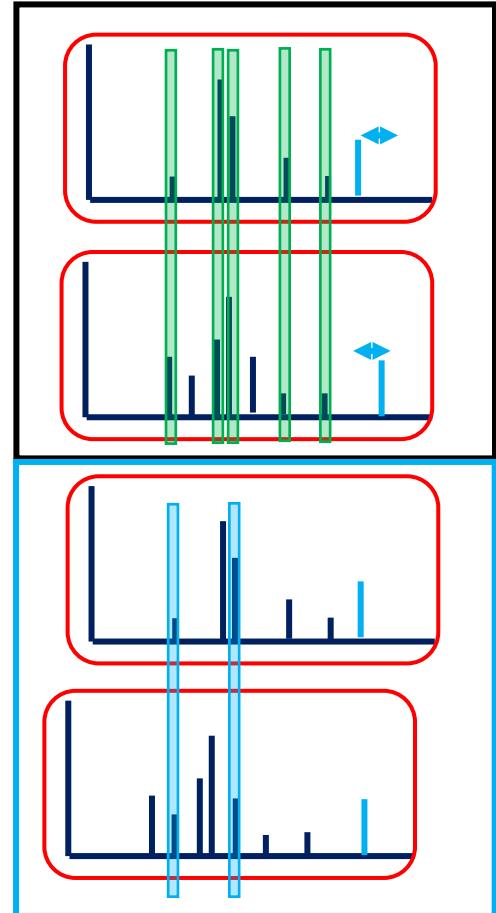
ClassyFire Scoring:

Flavonoids: Database hits:
2.25 nodes/ 4 nodes/
6 nodes = 6 nodes =

0.375 **0.67**

Illuminating the Rhamnaceae chemistry

Molecular Networking



plant related classifications:

different flavonoids

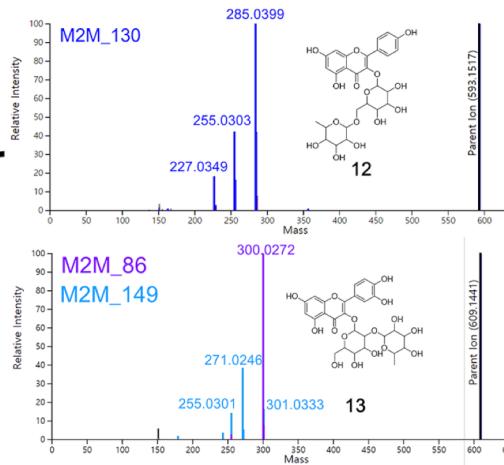
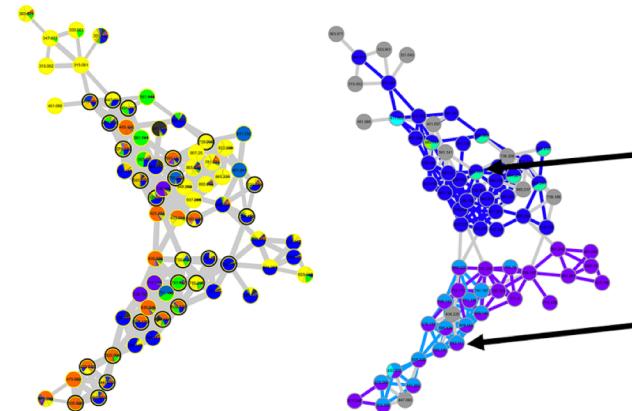
phenolic glycosides

triterpenoids

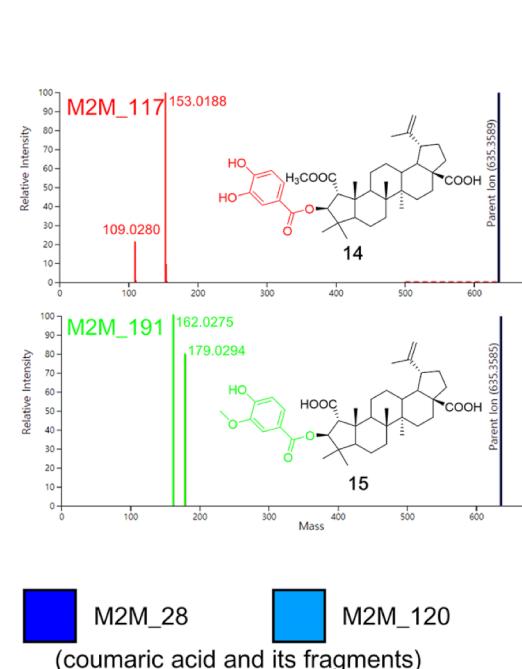
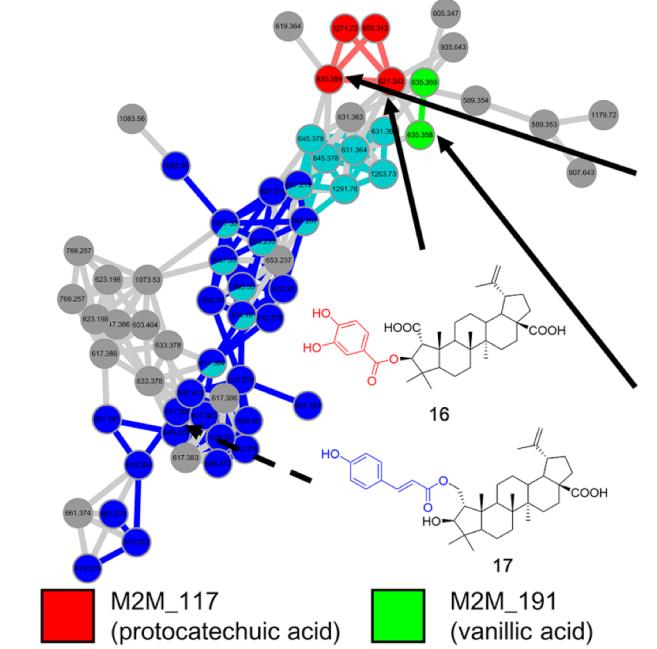
Dr Kyo Bin Kang, UCSD



Annotating plant molecular families



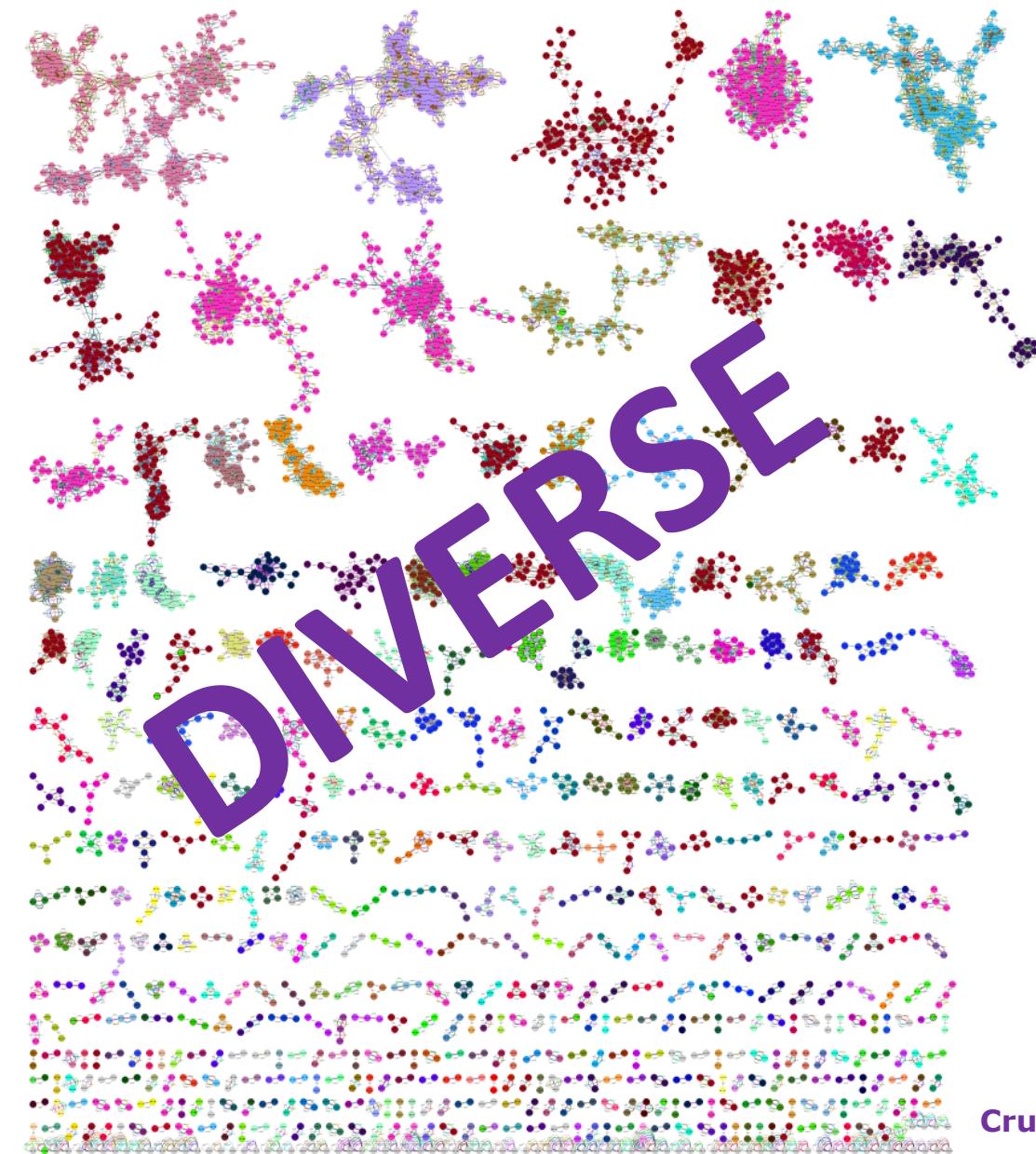
Flavonoid-3-O-glycoside Molecular Family:
Differentiation of subfamilies
Kaempferol and Quercetin based



Triterpenoid Molecular Family:
Differentiation of modifications
Protocatechuic acid and Vanillic acid based



Illuminating bacterial chemistry



DI
VERSE

**Molecular Network of
146 *Streptomyces* and *Salinispora* sp.**

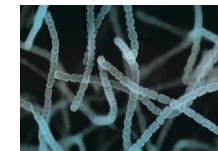
Triterpenoids

Lipids - PE

Cyclic peptides

Alpha amino acid esters

and many more....



Annotation of bacterial family I

Mass2Motif Annotation

You can assign a label (annotate) this Mass2Motif from the **Annotation** field below. Additionally, a shorter annotation can also be assigned through the **Short Annotation** field. This will be used in the network visualisation.

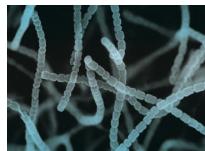
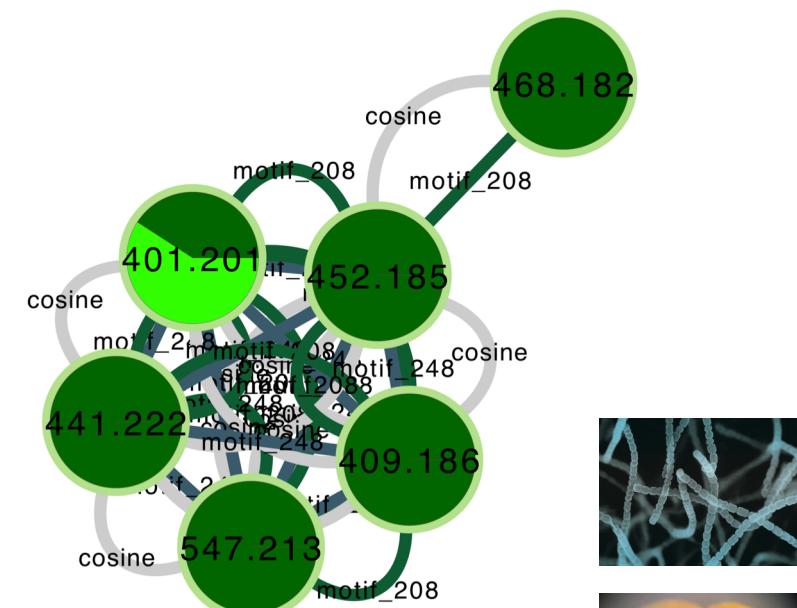
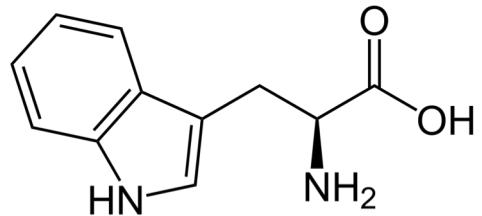
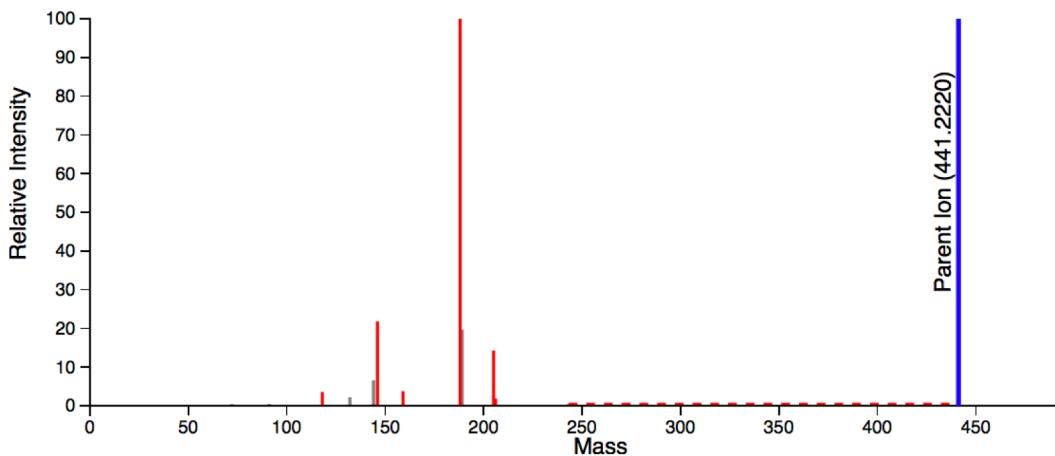
Annotation

Short Annotation

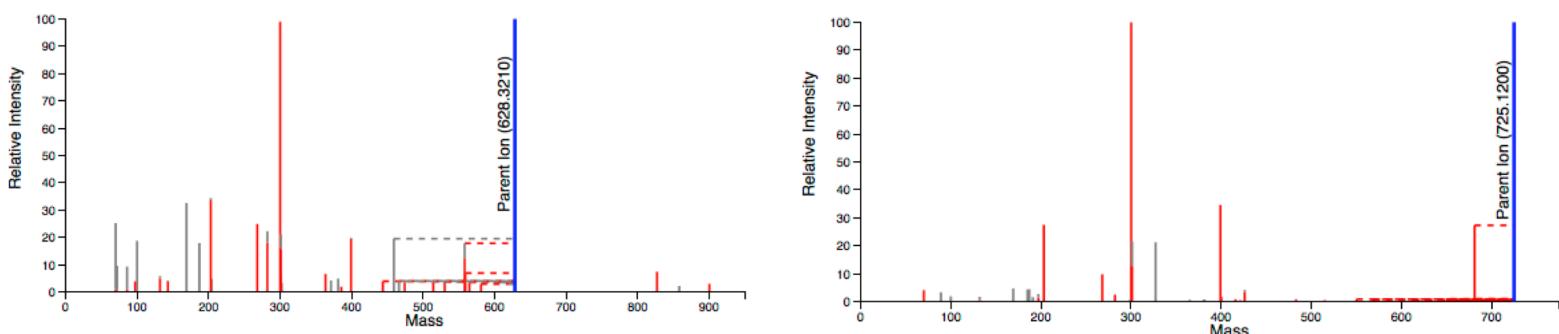
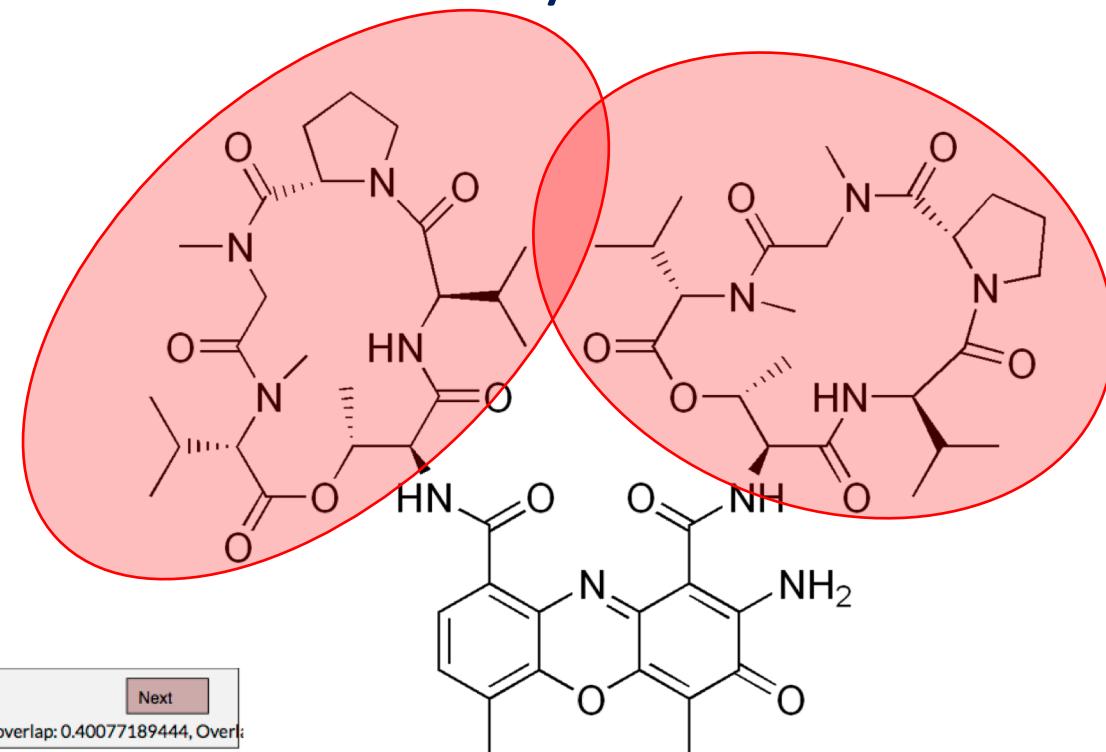
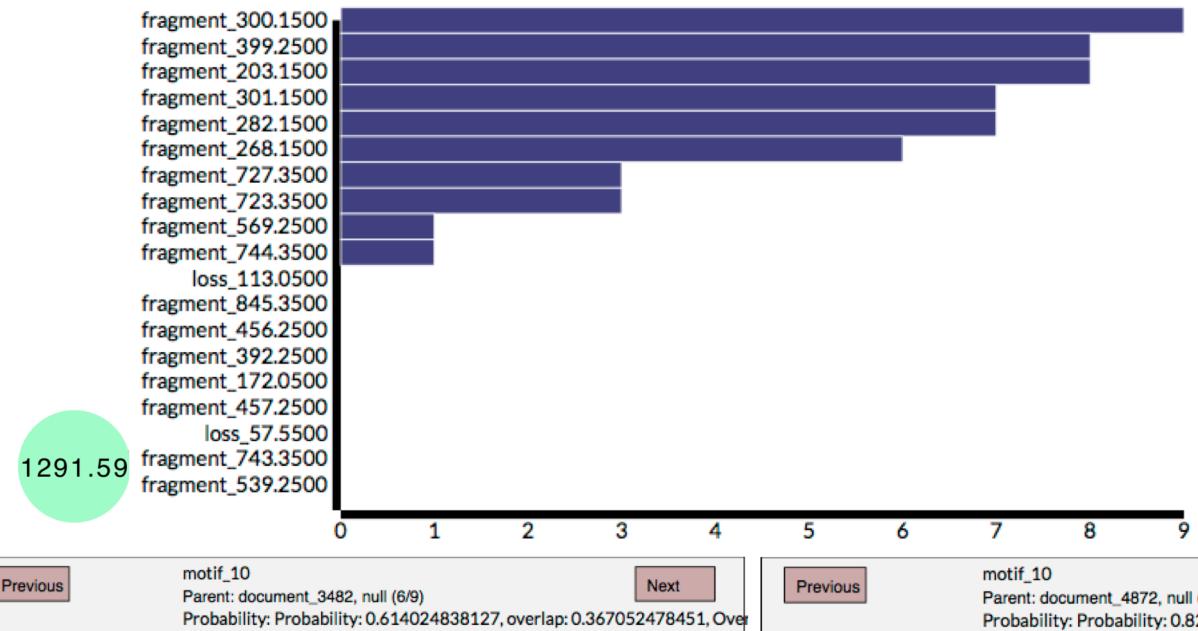
Save

Motif 208

Previous motif_208
Parent: document_1202, null (2/19)
Probability: Probability: 0.832888449166, overlap: 0.356384148972, Over

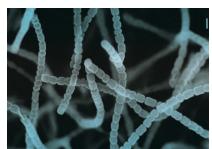


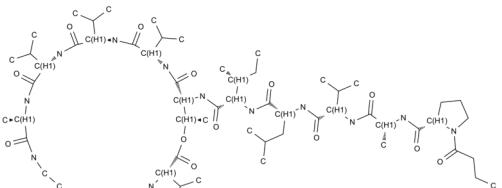
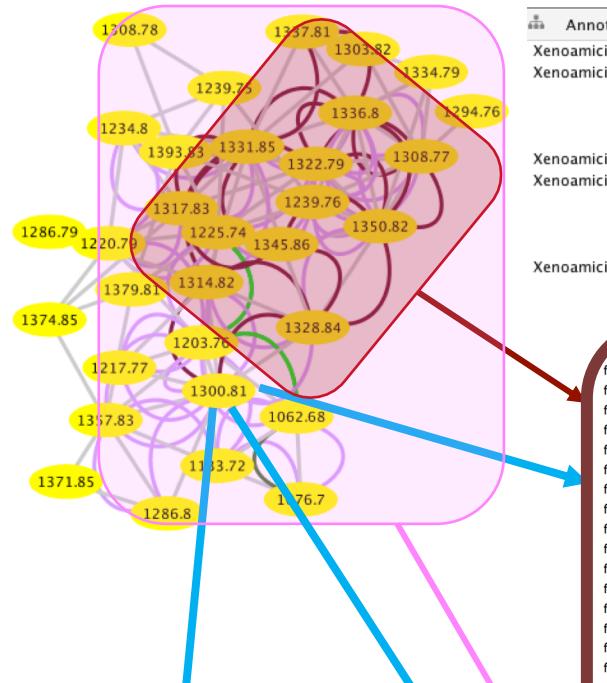
Annotation of bacterial molecular family II



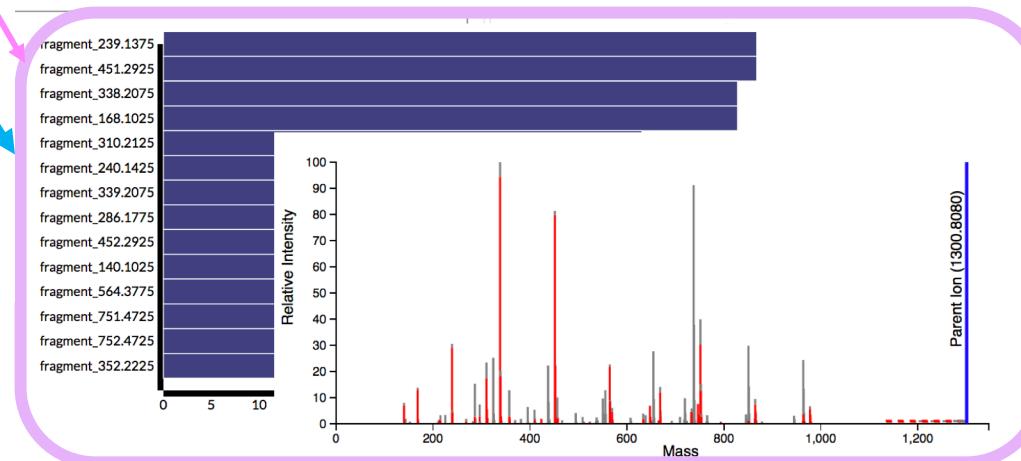
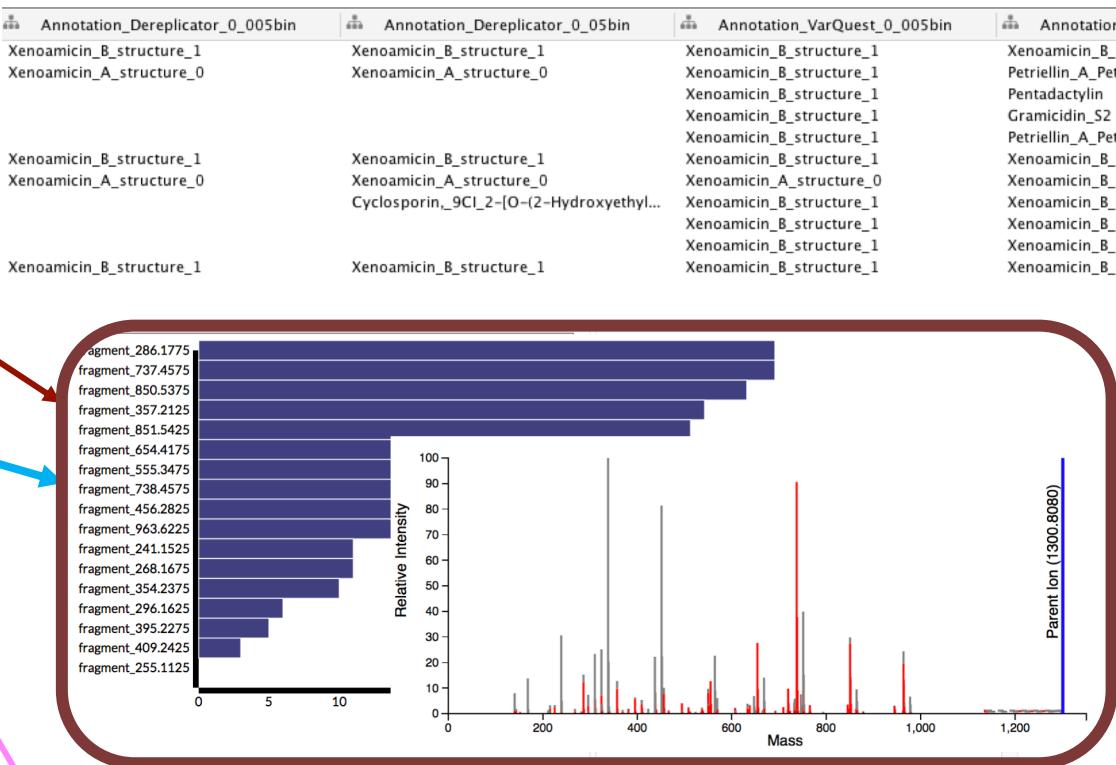
Actinomycin D

| PeptideMass | |
|-------------|---------|
| .3E-7 | 1414.84 |
| 7E-10 | 1268.57 |
| 1.0E-7 | 1290.57 |



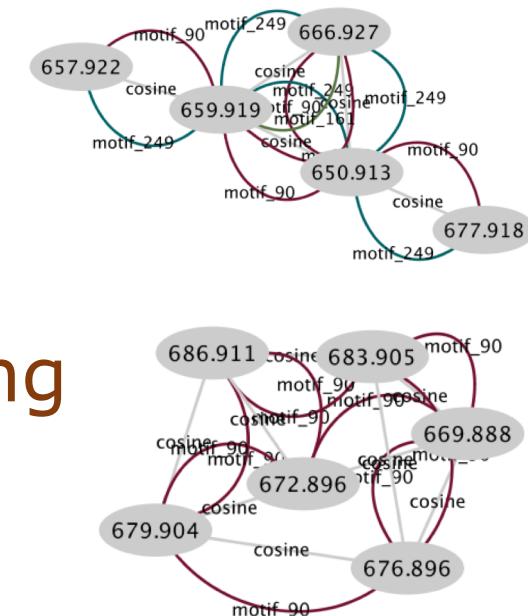


Xenoamicin A

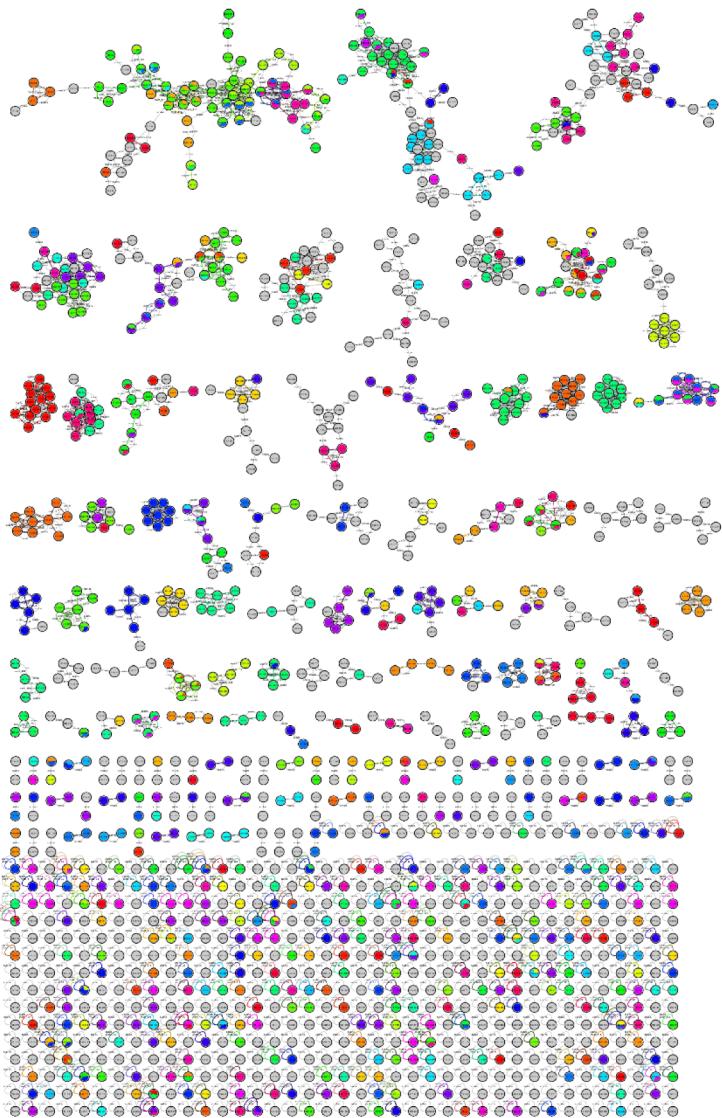


Ring

Tail



Annotation of fungal garden

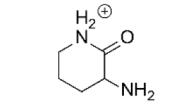


Unknown

COLOR and SHAPE CODE

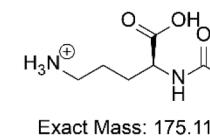
● Motif 1 - N2-acetylornithine related

Fragment 115.0875



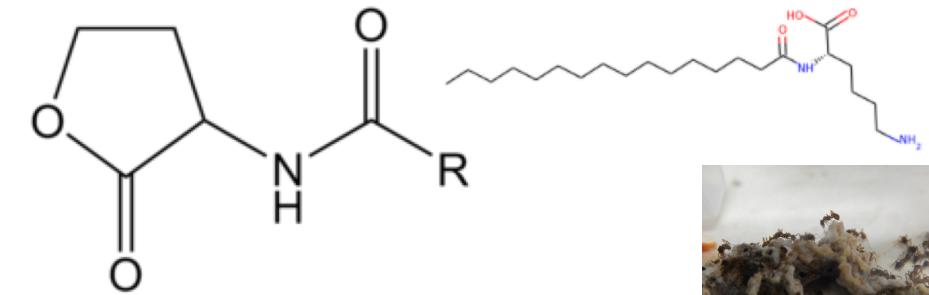
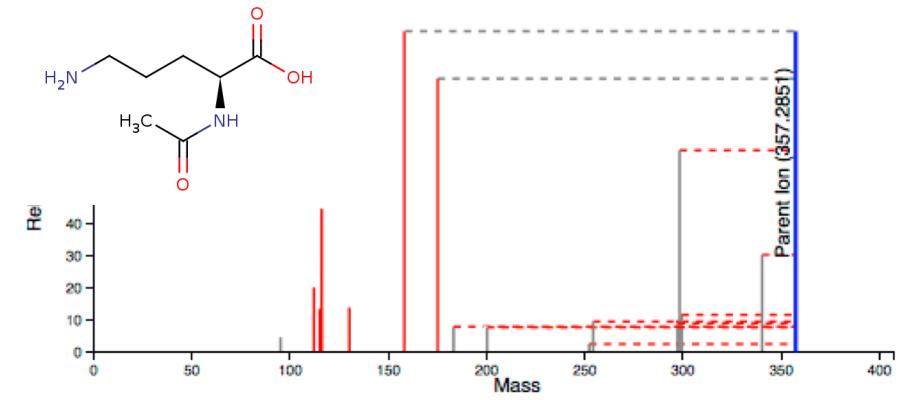
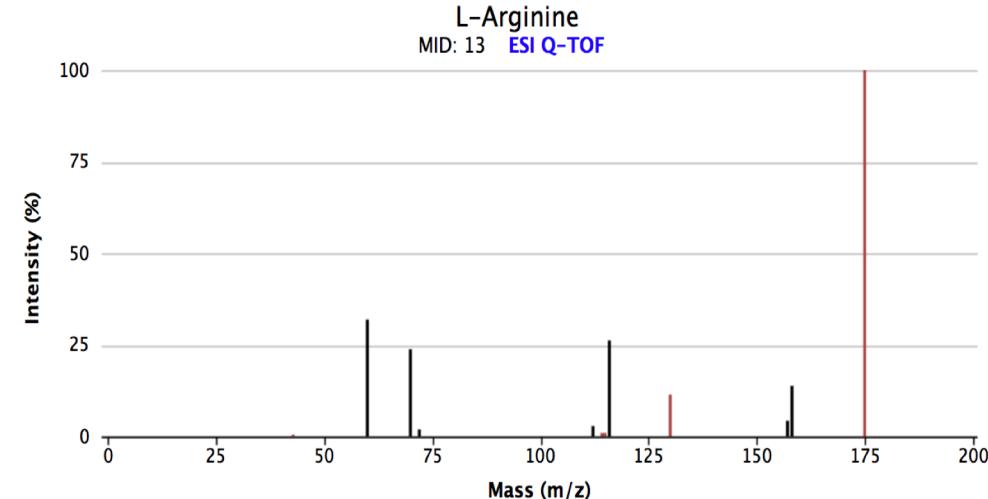
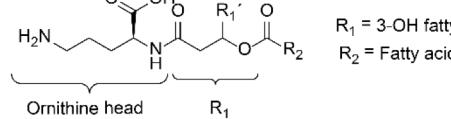
Exact Mass: 115.09

Fragment 175.1175

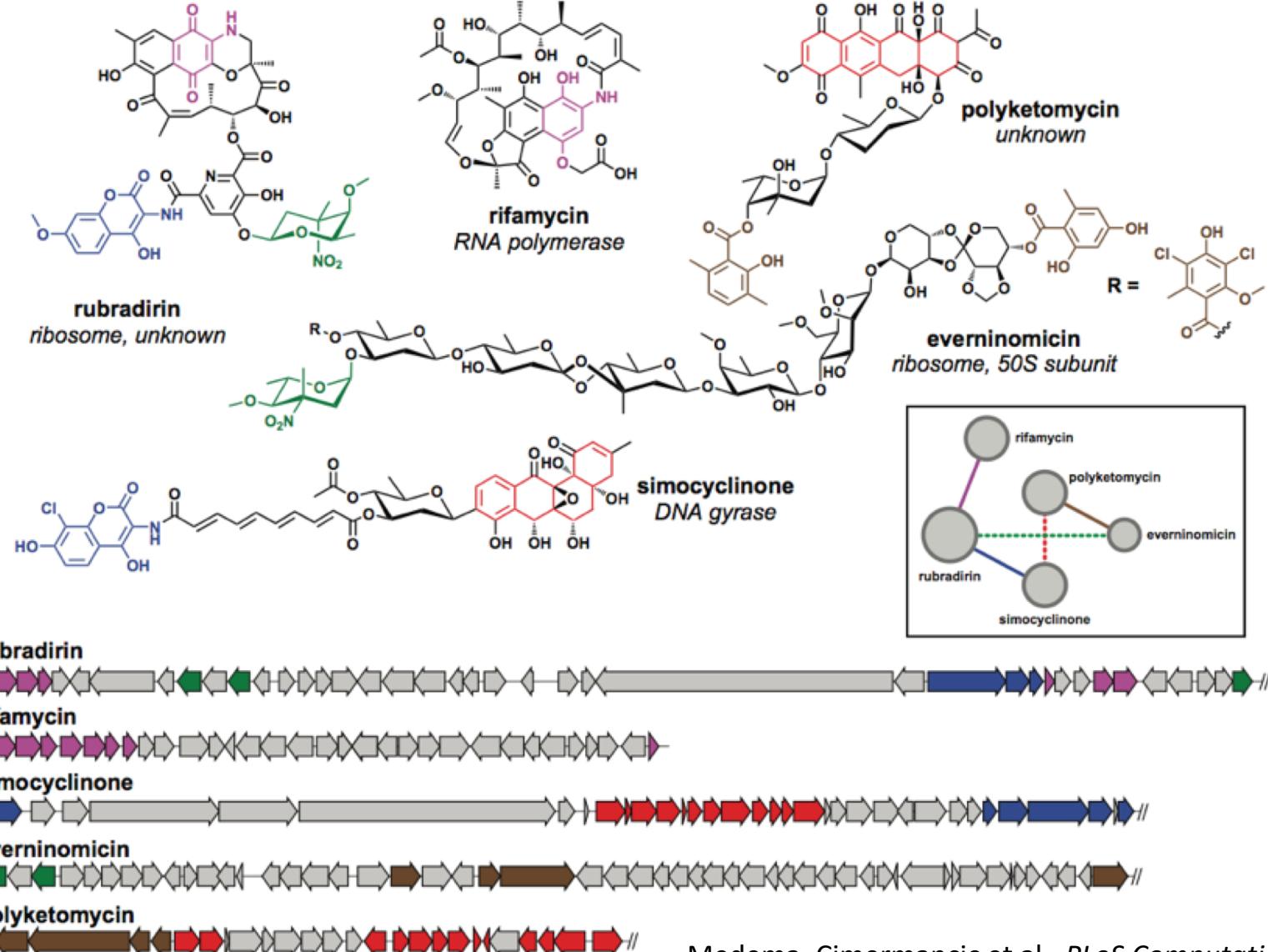


Exact Mass: 175.11

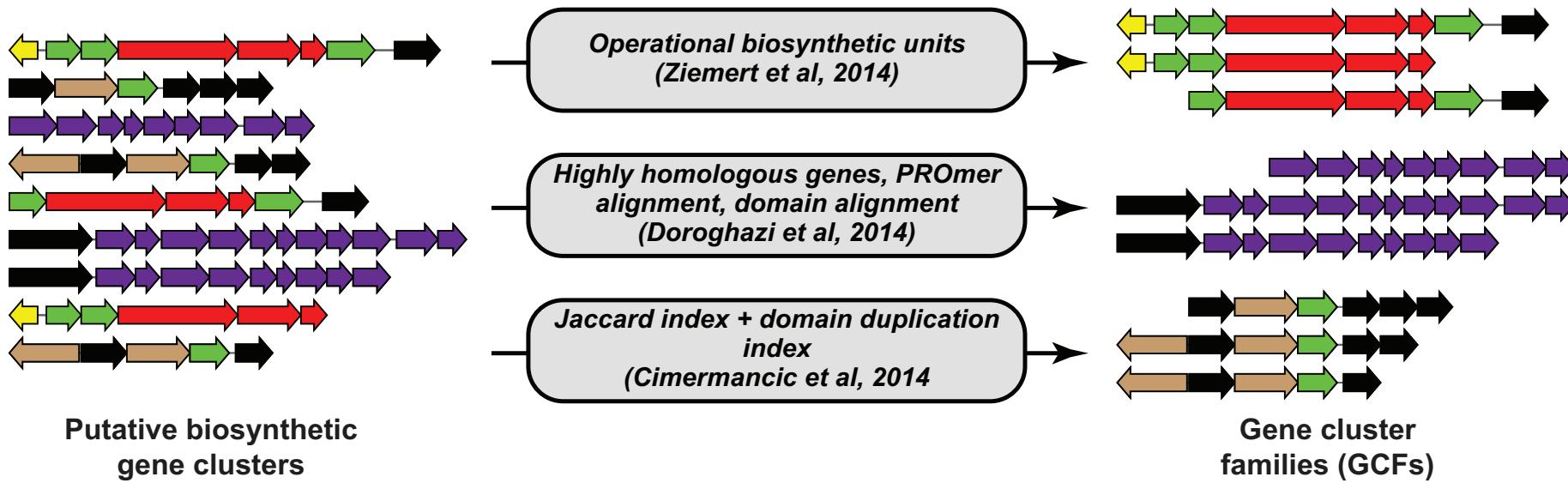
<https://metlin.scripps.edu>



Biosynthetic gene clusters: the key to mine genomes for novel natural product chemistry

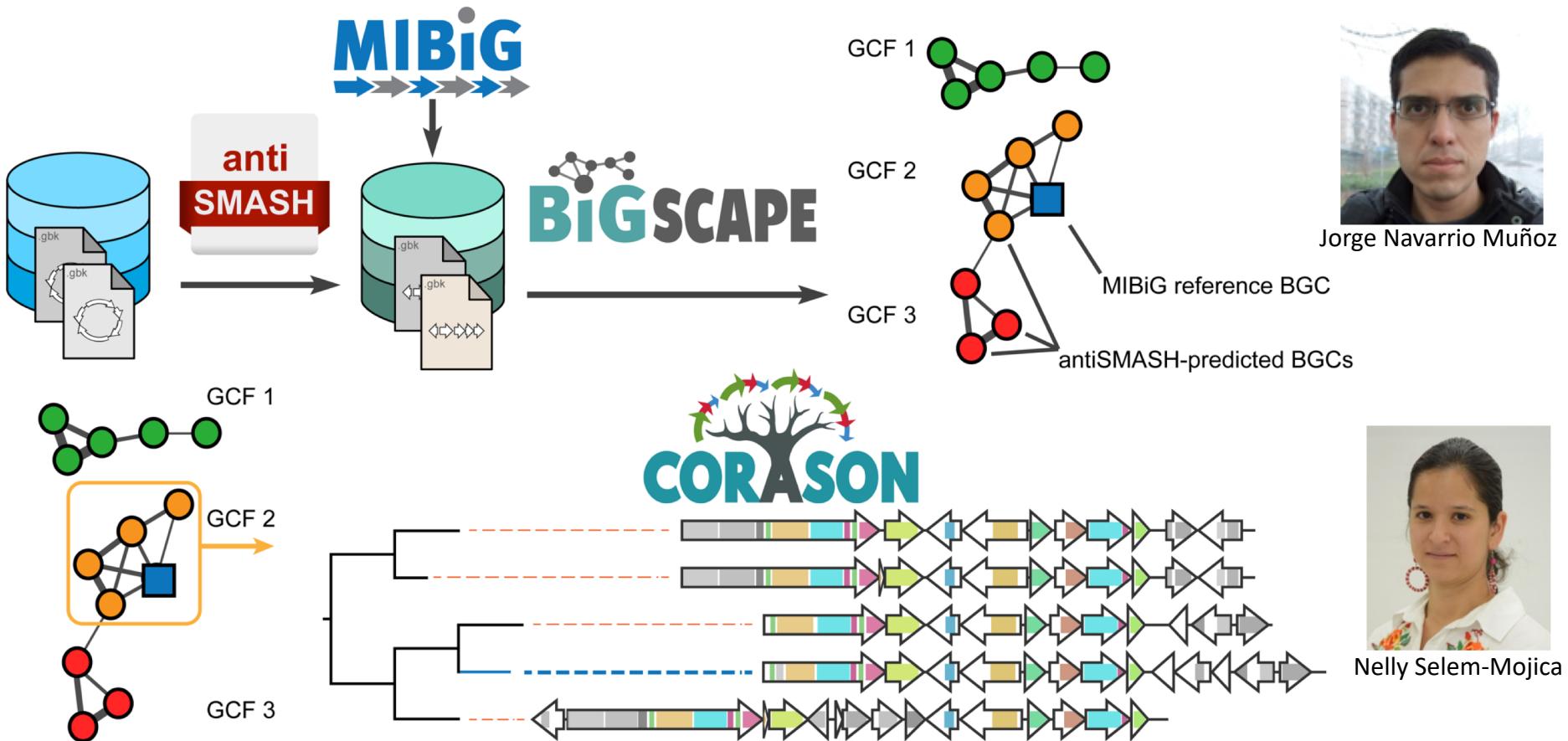


Grouping Biosynthetic Gene Clusters into Families

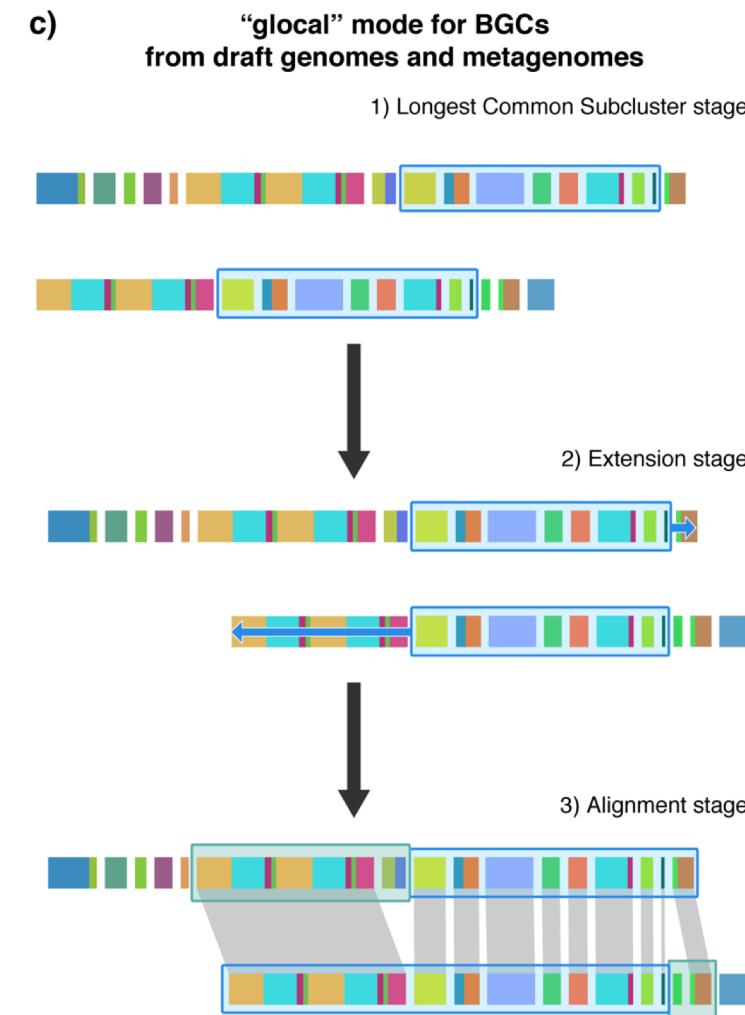
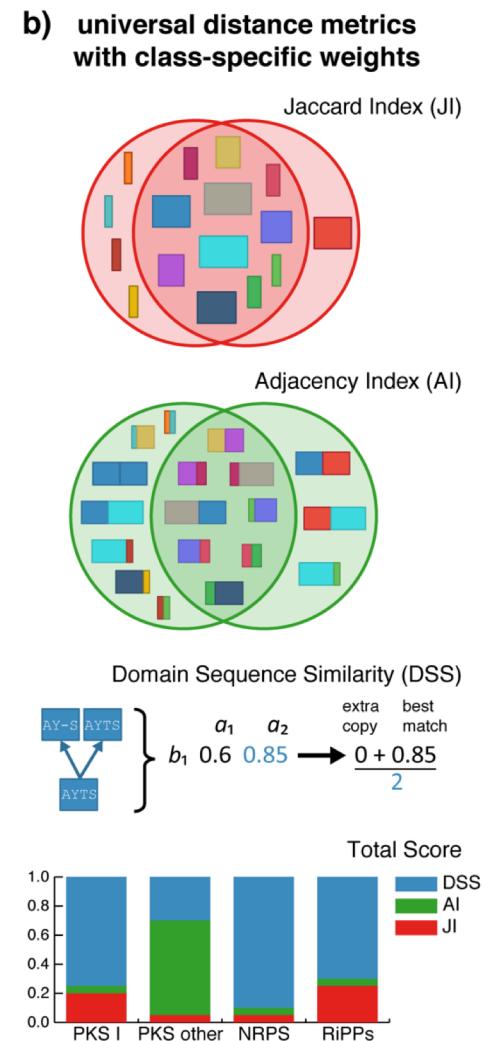
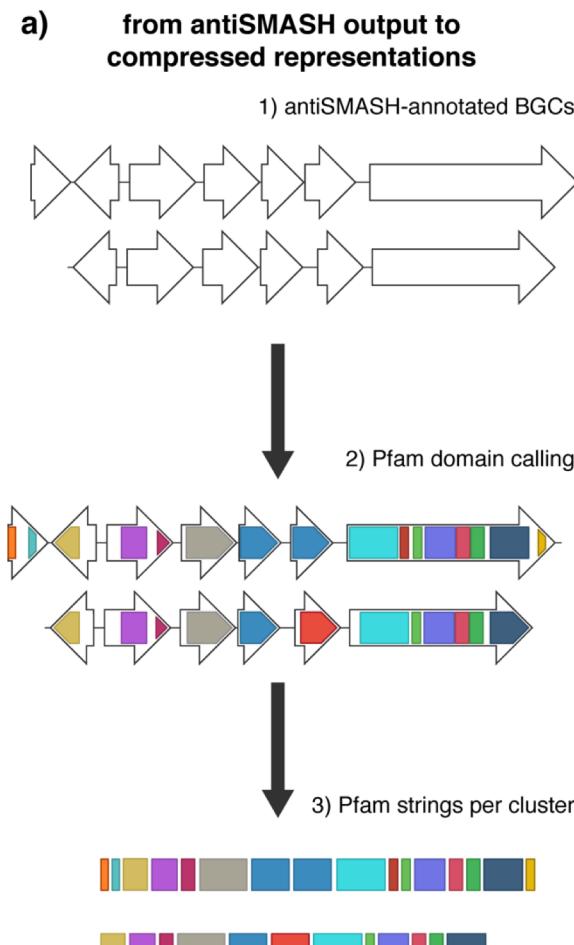


24

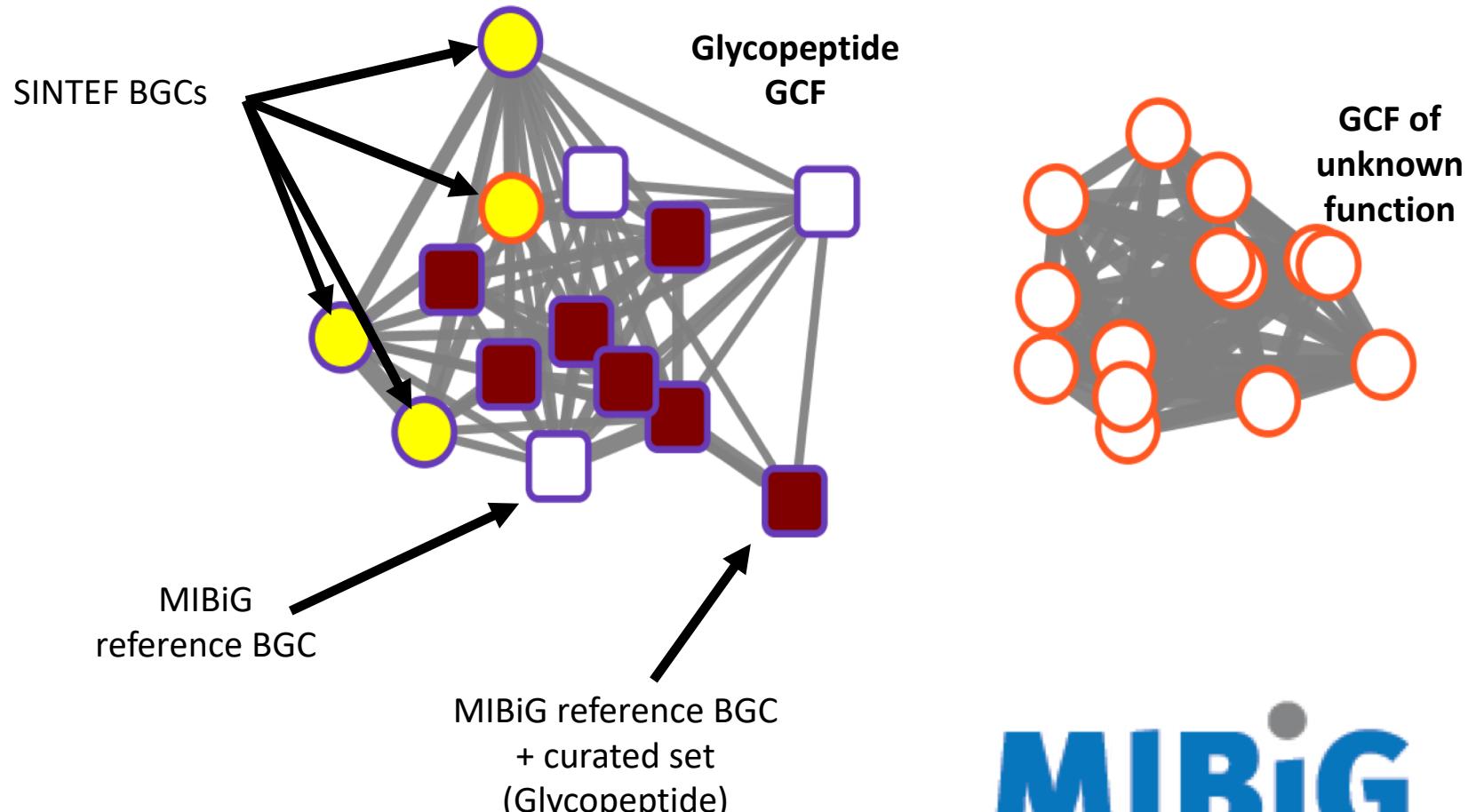
BiG-SCAPE & CORASON automate reconstruction and phylogenomic analysis of BGC families



Efficient Data compression, Distance metrics and alignment modes

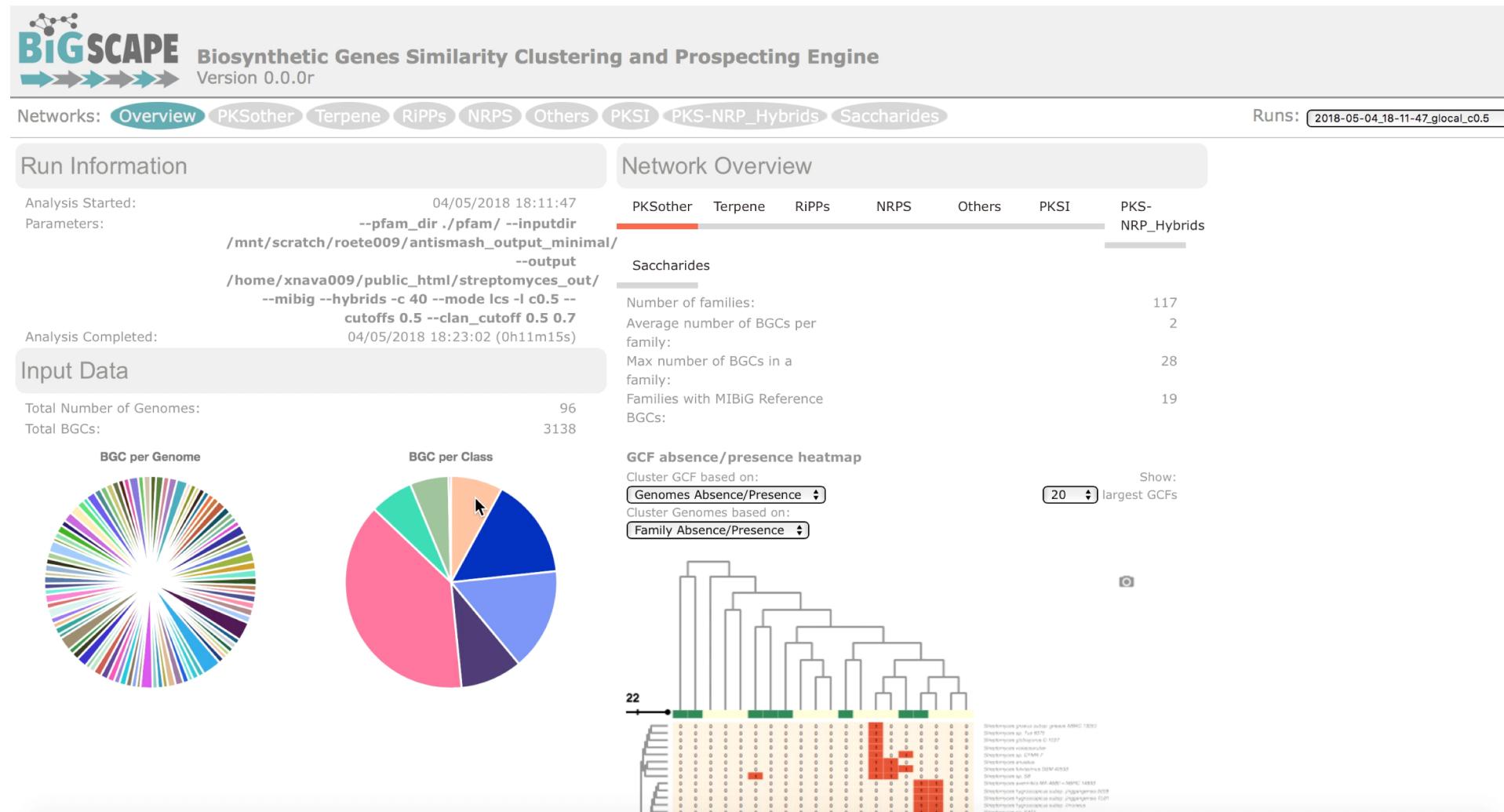


MIBig reference data allows rapid annotation propagation



<https://mibig.secondarymetabolites.org>

A rich web interface allows effectively exploring biG-scape outputs



WAGENINGEN
UNIVERSITY



Satria Kautsar

<https://git.wageningenur.nl/medema-group/BiG-SCAPE>

Big-scape has quickly been adopted by the scientific community

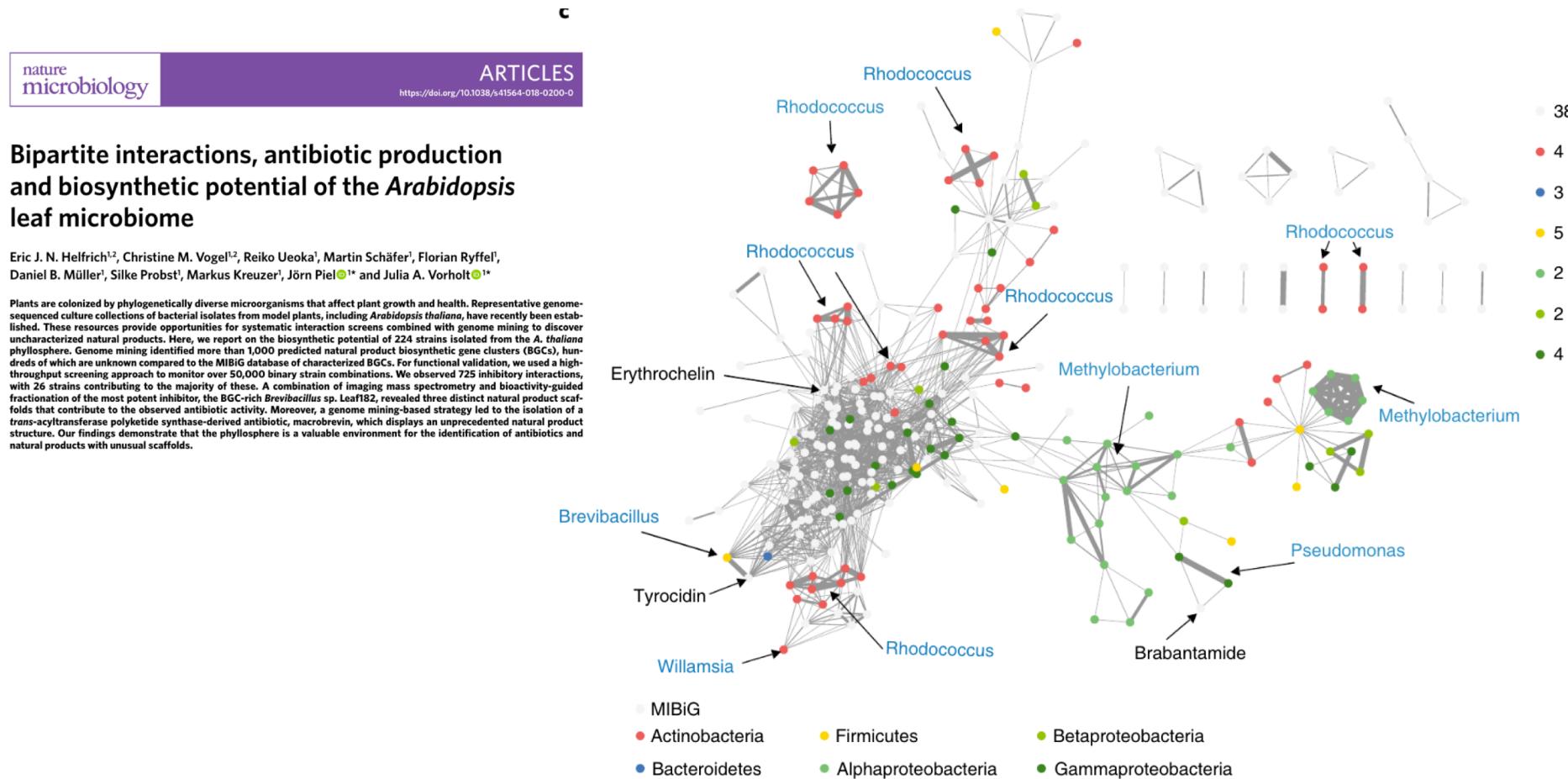
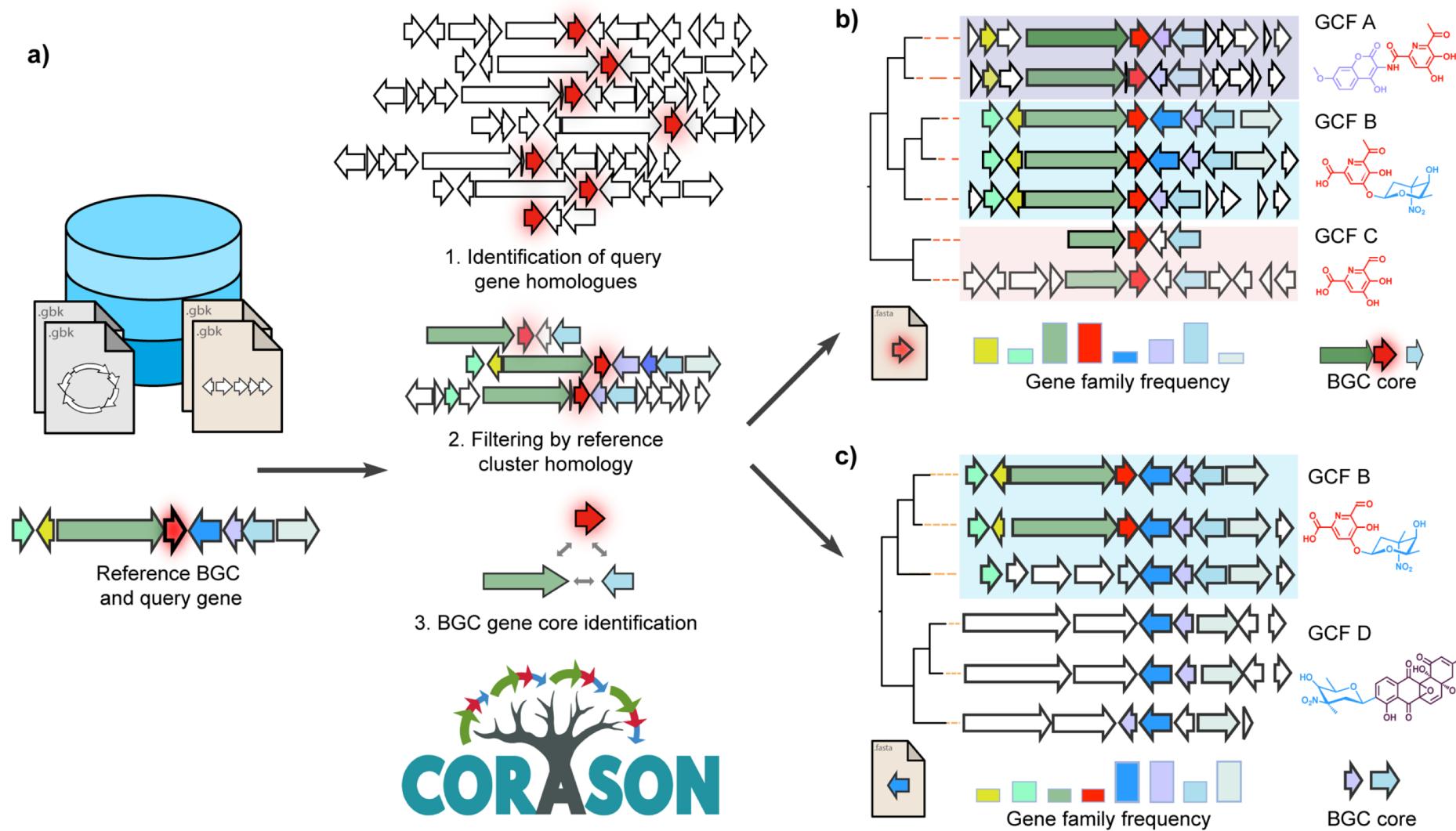
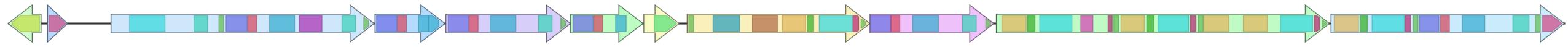


Fig. 3 | BiG-SCAPE analysis of BGCs detected by antiSMASH in 207 genomes of the At-LSPHERE strain collection and comparison with the MIBiG database of characterized BGCs. a, RiPP BGCs ($n=246$). b, Type I PKS BGCs ($n=277$). c, NRPS BGCs ($n=338$). Nodes and singletons are colour coded according to their origin (MIBiG database or phylum/class) and represent individual BGCs. The widths of interconnecting edges indicate the degree of relatedness between two BGCs, with connections up to a raw distance of 0.75 retained (Supplementary Table 12). Numbers of singletons are indicated. Black labels denote compounds associated with selected MIBiG BGCs. Blue labels highlight clusters of related BGCs or individual BGCs of interest by phylogenetic distribution.

CORASON facilitates finding gene clusters around homologues of gene of interest



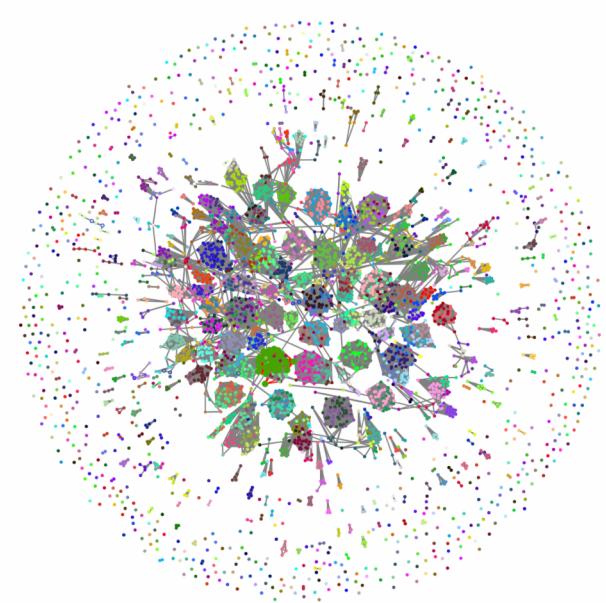
Linking substructures to genetic elements



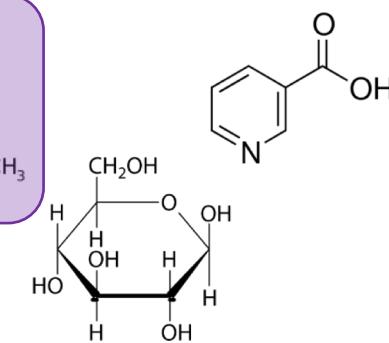
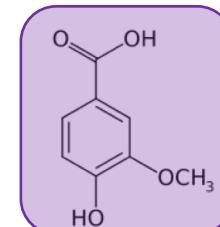
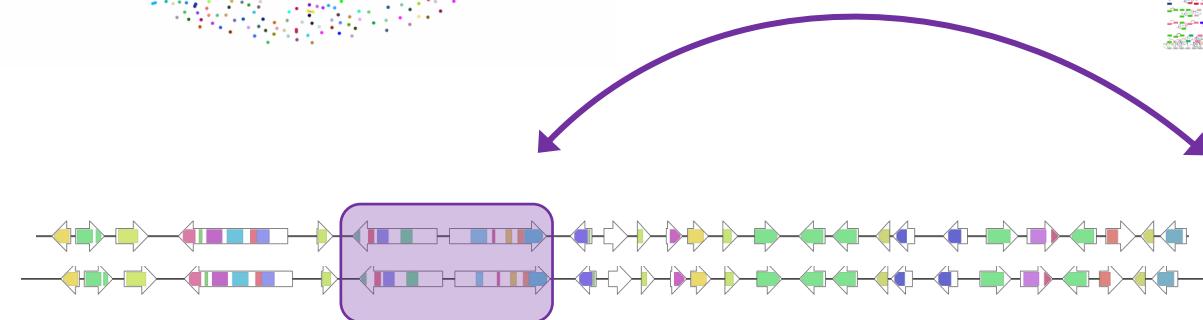
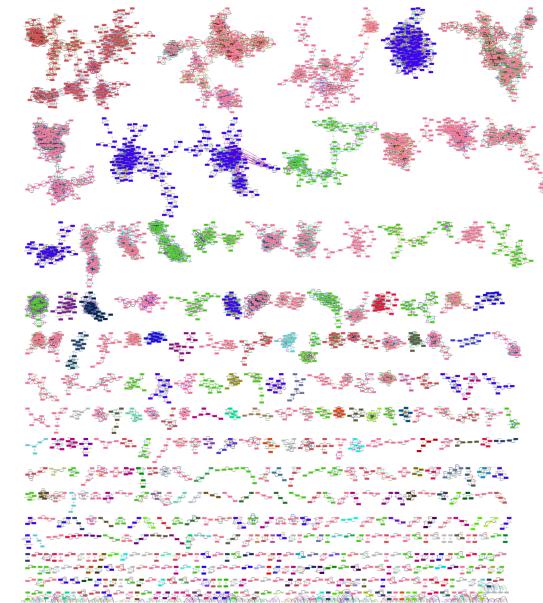
iOMEGA: Integrated Omics for MEtabolomics and Genomics Annotation

■ Gene Cluster Families &

Metabolite Families



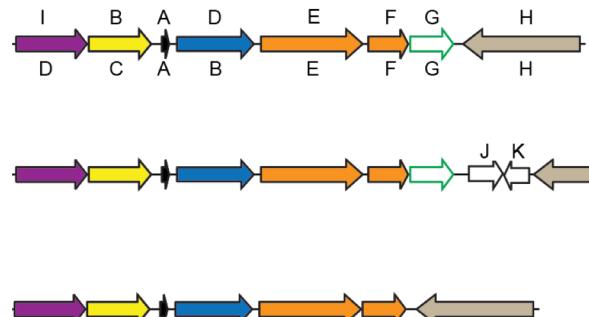
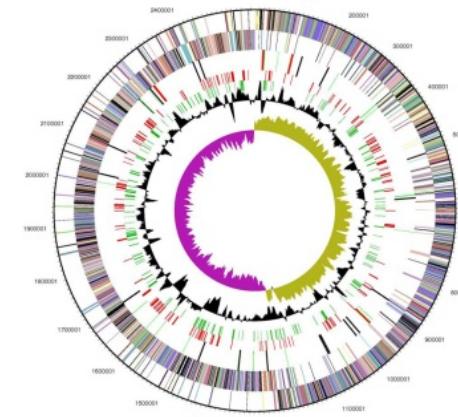
<-->



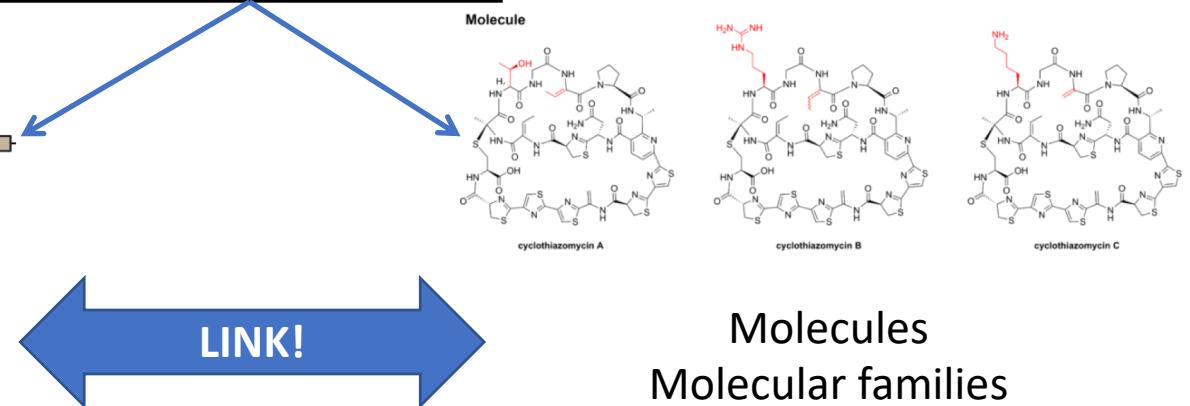
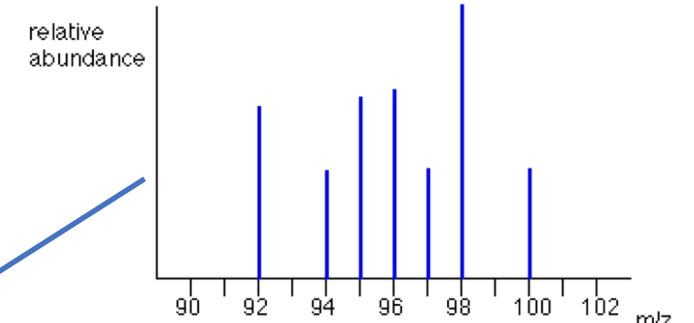
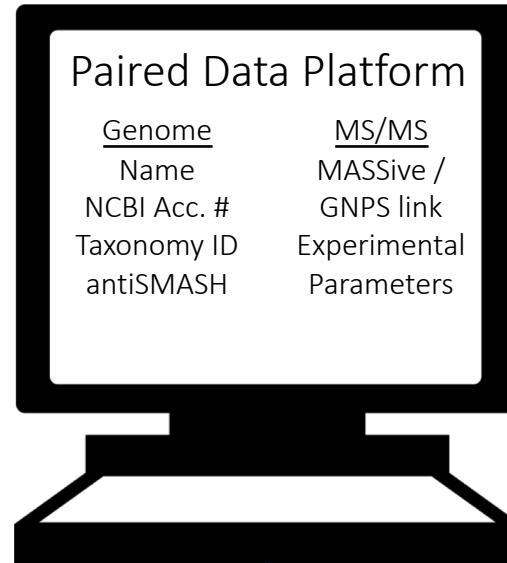
Paired Data Platform



Michelle Schorn



Gene clusters
Gene cluster families



Paired Data Platform

iOMEGA paired data platform schema JuNo

This is the JSON schema for paired genomic / metabolomic data.

version *

1

Personal data*

Name of contact for correspondence

This person will be the point of contact for any communication related to this entry.

Academic institution or company name

Please use the full, official name of your institute in English. E.g., 'Harvard University'.

Submitter contact e-mail address

Name of the principal investigator of the submitter

This person is contacted in case the submitter has moved institution

Paired Data Platform

Extraction solvent



Please select the organic solvent used to extract the sample. If your solvent is not listed, please choose Polar or Non-polar. If you used multiple solvents, please select and order them and indicate the ratio below.

Methanol

Methylene Chloride / Dichloromethane

Ethyl acetate

Chloroform

Acetone

Isopropanol

Butanol

Acetonitrile

+ ratio here.



Thank you!

Special Issue: Metabolomics Data Processing and Data Analysis—Current Best Practices

Deadline for manuscript submissions: **28 February 2019**

Guest Editors



Dr. Kati Hanhineva
University of Eastern
Finland
Twitter: @KatiHanhineva



Dr. Justin van der Hooft
Wageningen University
Netherlands
Twitter: @jjvanderhooft

Keywords

Metabolomics data processing data interpretation annotation and visualization data analysis