

# A community-driven paired data platform to accelerate natural product mining

**Justin J.J. van der Hooft et al.**

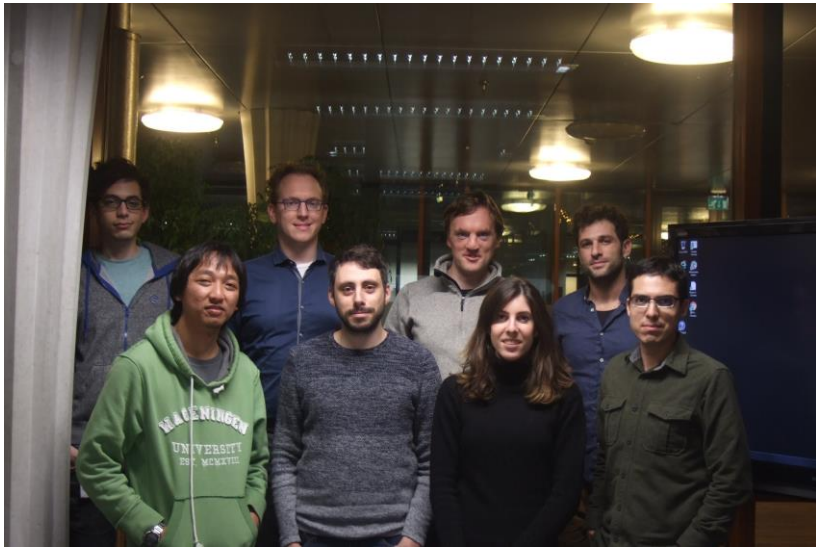
Bioinformatics Group – Wageningen University, The Netherlands

Wageningen, 12 December 2018





# Team work! 😊



**Medema lab - Wageningen UR, NL**



**NL eScience Center**



**Dorrestein lab - San Diego, USA**

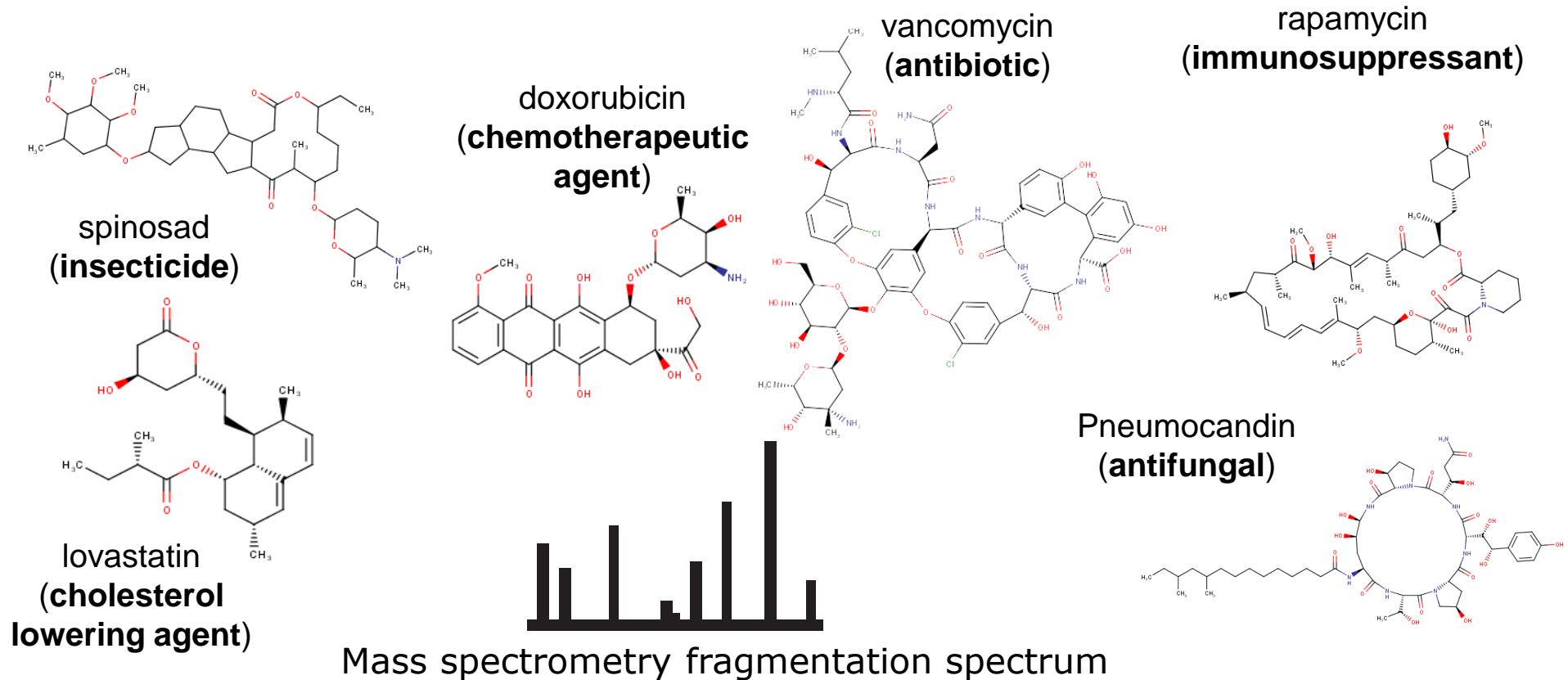


**Glasgow Polyomics - University of Glasgow, UK**



# The challenge in metabolomics....

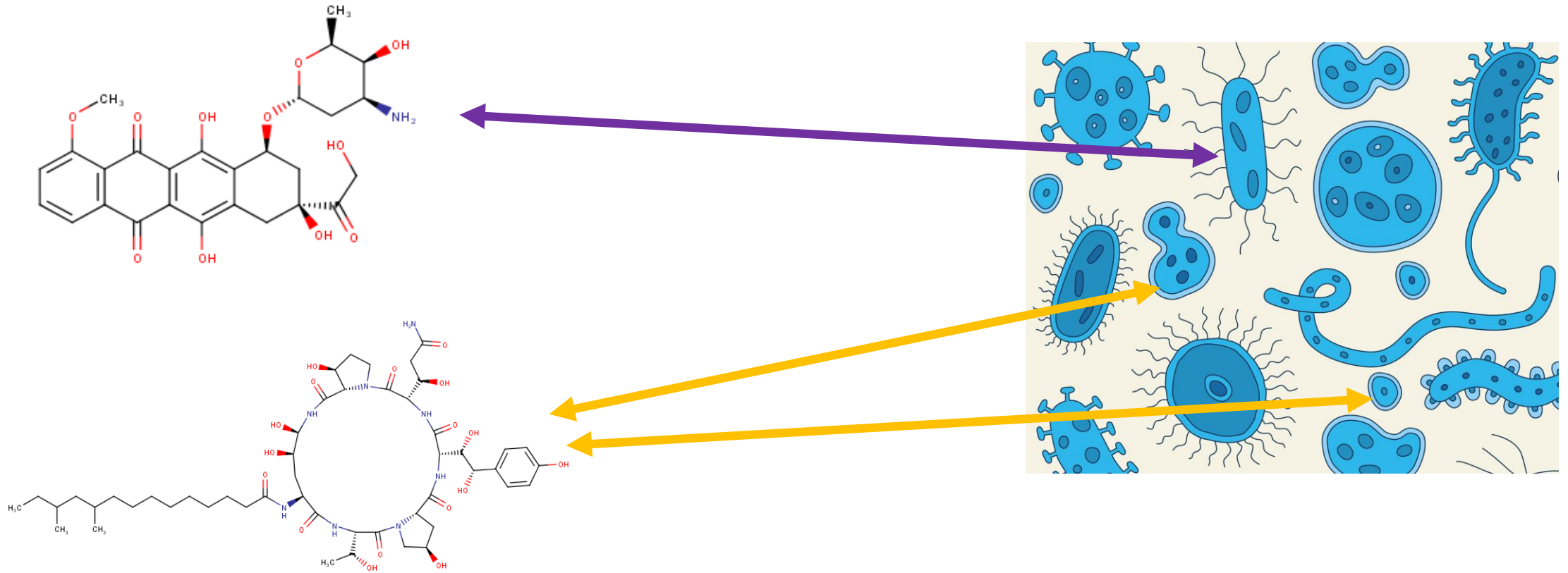
**Nature produces a large & diverse arsenal of high-value molecules:**



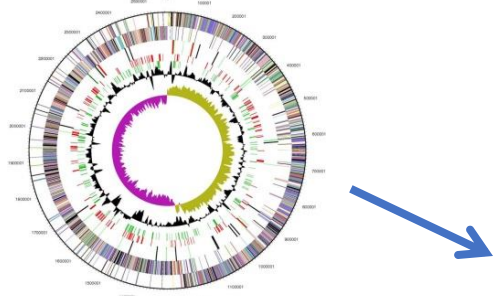
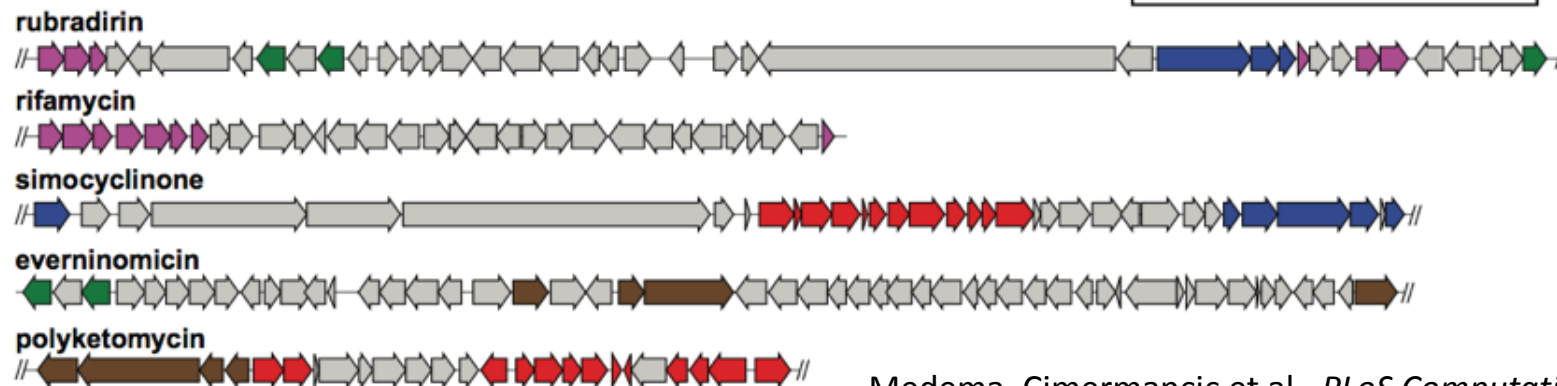
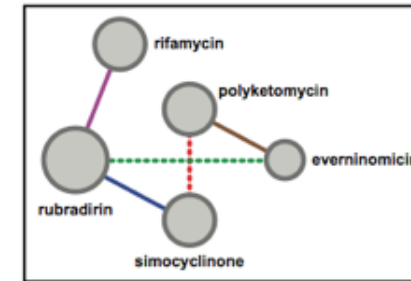
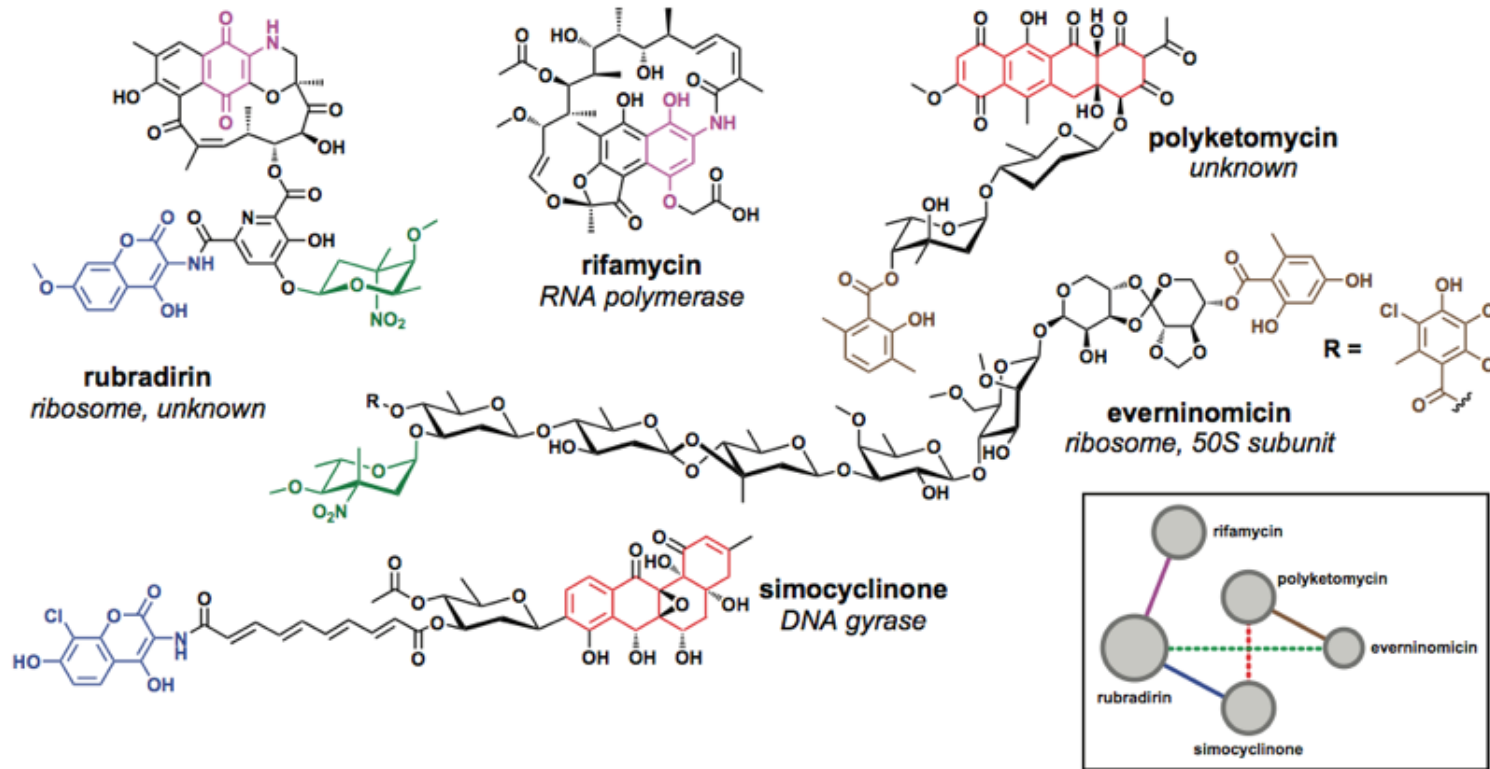
....is **large-scale coupling of spectral data to molecular structures**  
of known & especially **novel** natural products molecules



# iOMEGA project: integrated omics for metabolomics and genomics annotation



# Biosynthetic gene clusters: key to mining genomes for chemistry



State of the art....

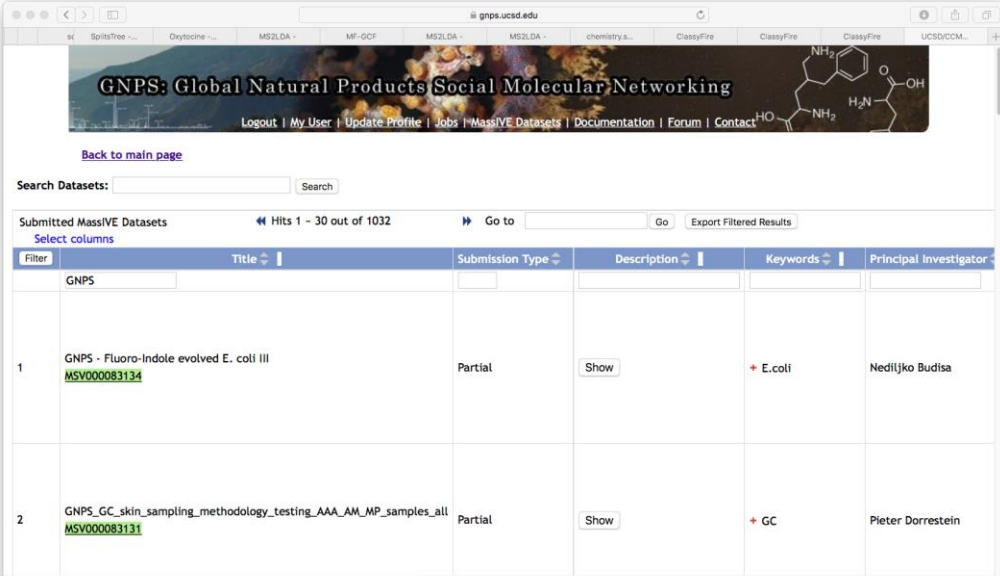
Public genome repositories

Public metabolome repositories

Libraries of validated biosynthetic gene clusters

Libraries of annotated and identified molecular spectra

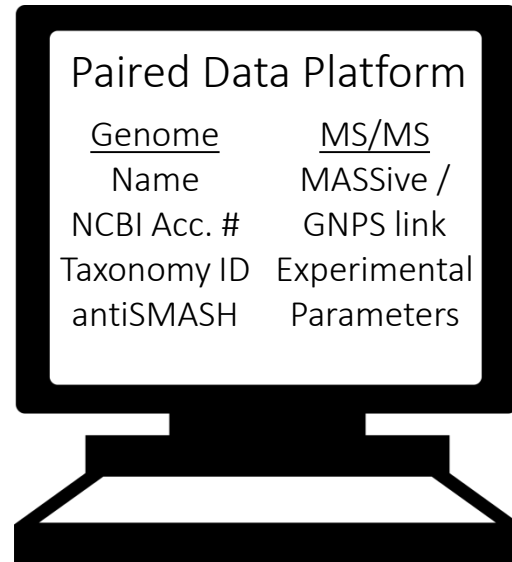
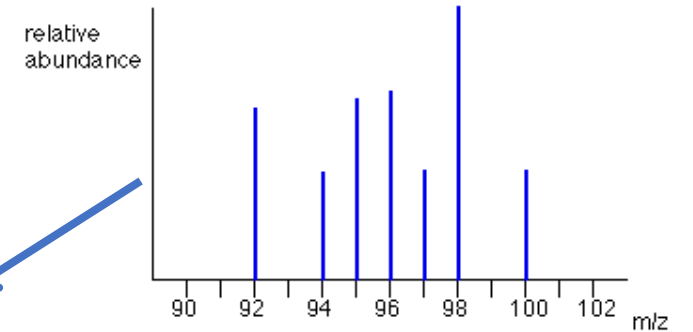
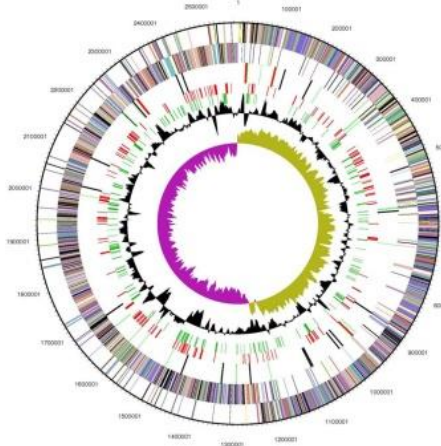
How to link it all?



The screenshot shows the GNPS (Global Natural Products Social Molecular Networking) website interface. At the top, there is a navigation bar with the site name and a chemical structure. Below the navigation bar, there is a search bar and a table of submitted datasets. The table has columns for Filter, Title, Submission Type, Description, Keywords, and Principal Investigator. Two datasets are visible in the table:

Filter	Title	Submission Type	Description	Keywords	Principal Investigator
GNPS	GNPS - Fluoro-Indole evolved E. coli III <a href="#">MSV000083134</a>	Partial	Show	+ E.coli	Nediljko Budisa
	GNPS_GC_skin_sampling_methodology_testing_AAA_AM_MP_samples_all <a href="#">MSV000083131</a>	Partial	Show	+ GC	Pieter Dorrestein

# Paired Data Platform

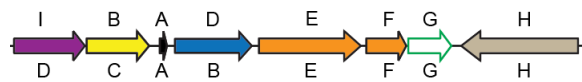


Paired Data Platform

<u>Genome</u>	<u>MS/MS</u>
Name	MASSive /
NCBI Acc. #	GNPS link
Taxonomy ID	Experimental
antiSMASH	Parameters

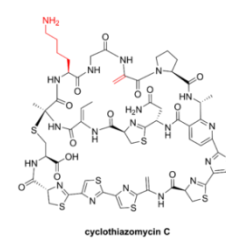
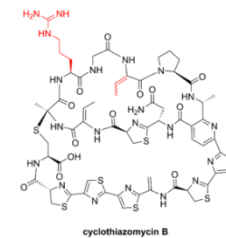
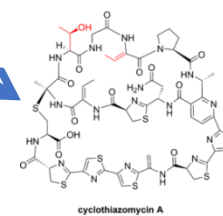


Dr. Michelle Schorn, WUR



Gene clusters  
Gene cluster families

Molecule



Molecules  
Molecular families



# Paired Data Platform

## iOMEGA paired data platform schema JuNo

---

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

**version \***

## 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.

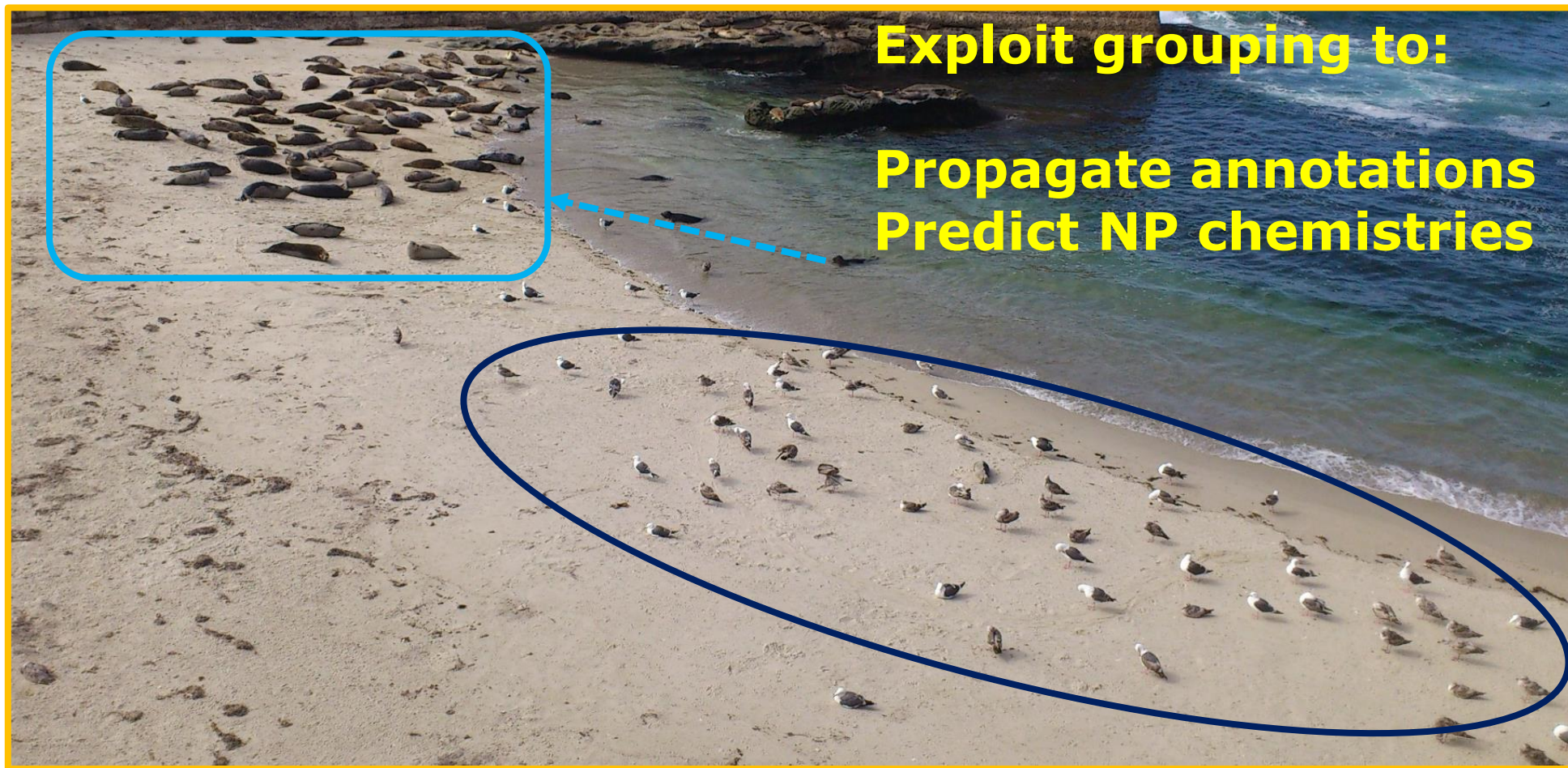


- Methanol**
- Methylene Chloride / Dichloromethane
- Ethyl acetate
- Chloroform
- Acetone
- Isopropanol
- Butanol
- Acetonitrile



ratio here.

# Improved annotation power by pattern mining and networking



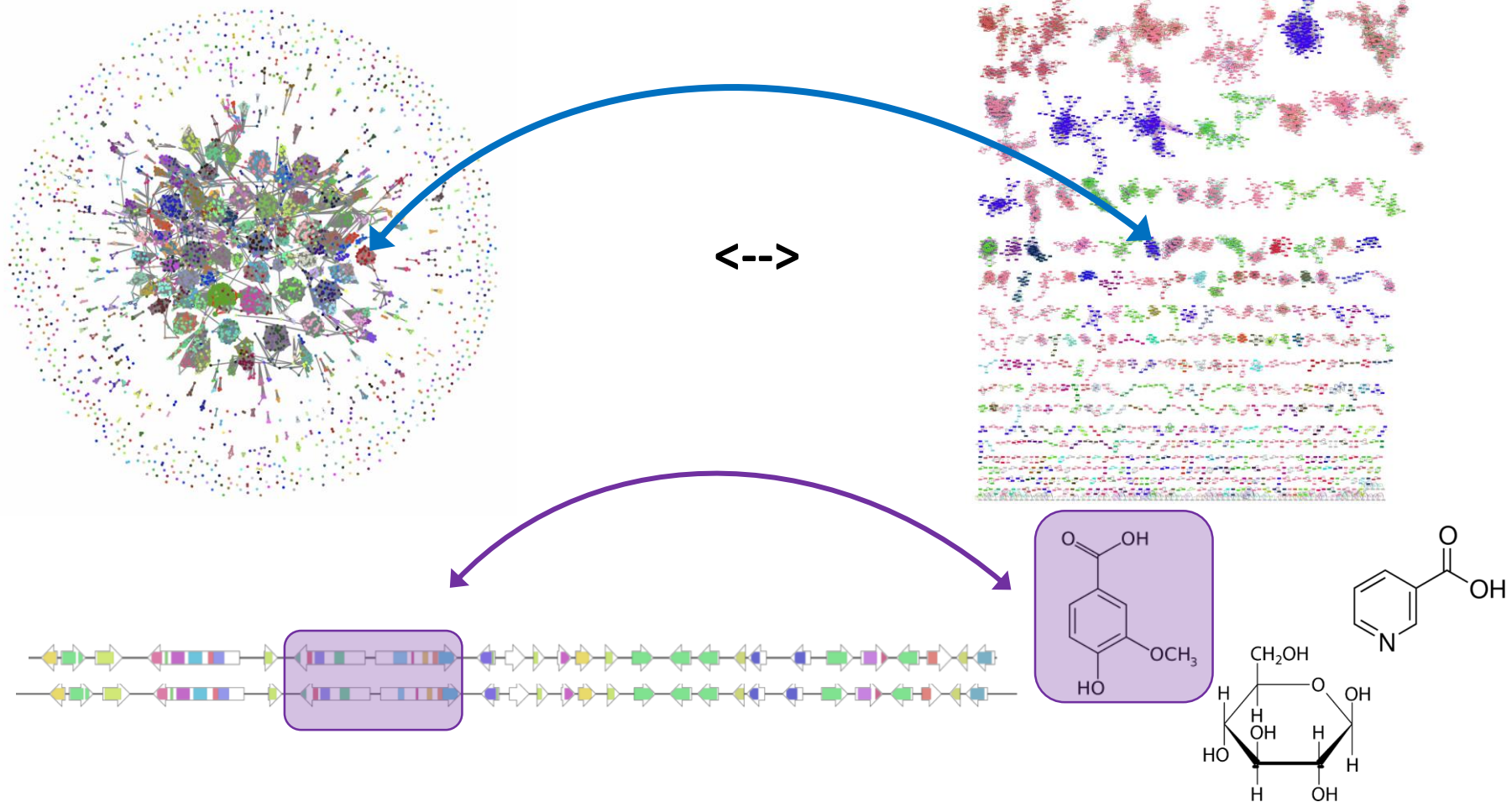
# Linking substructures to genetic elements

iOMEGA: Integrated Omics for MEtabolomics and Genomics Annotation

Gene Cluster Families

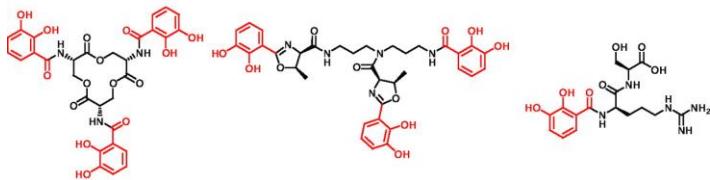
&

Metabolite Families

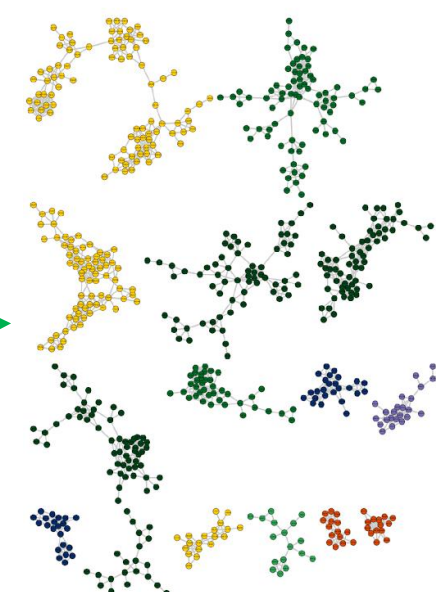
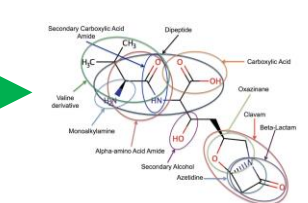
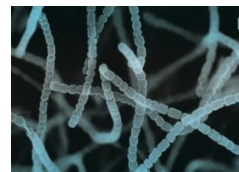
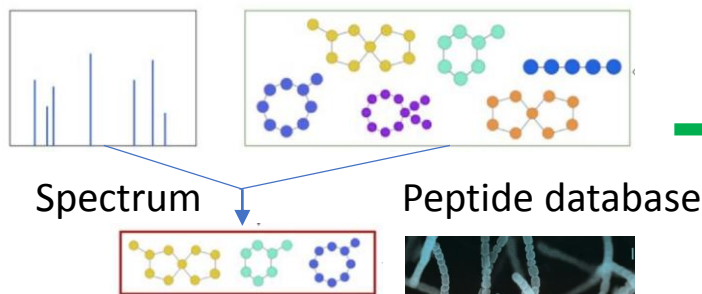
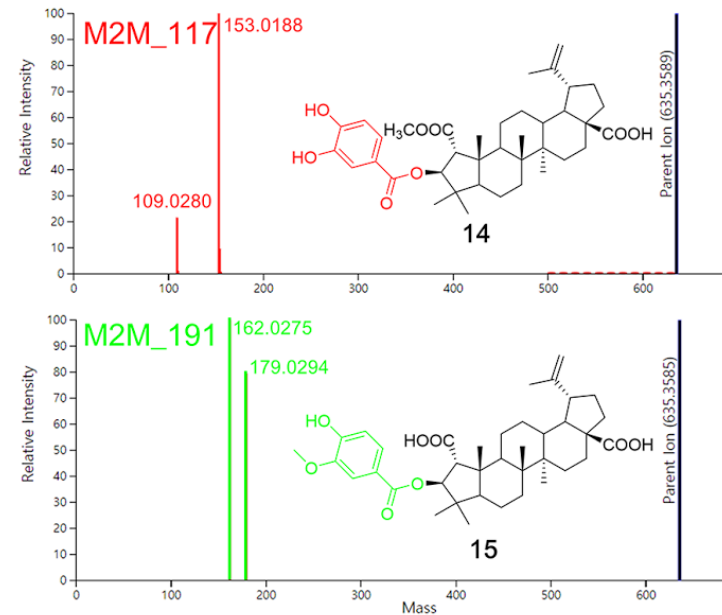
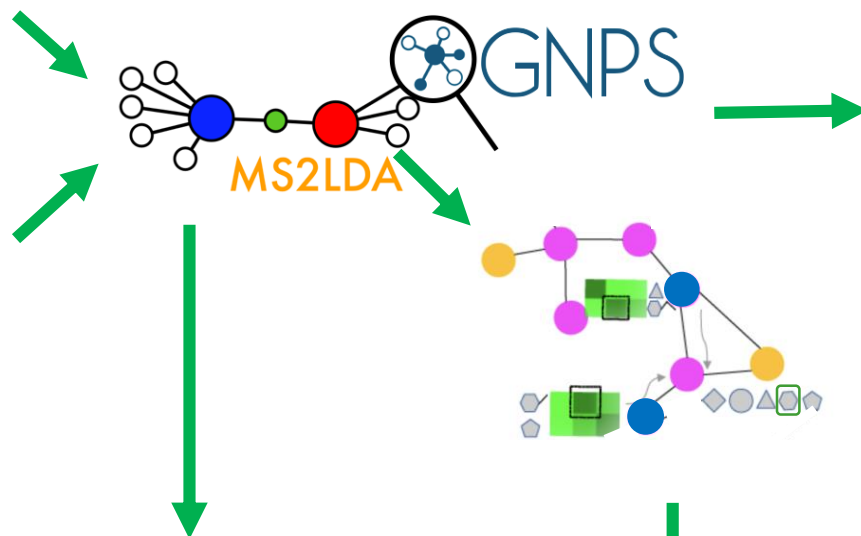




# Integrated metabolomics workflow



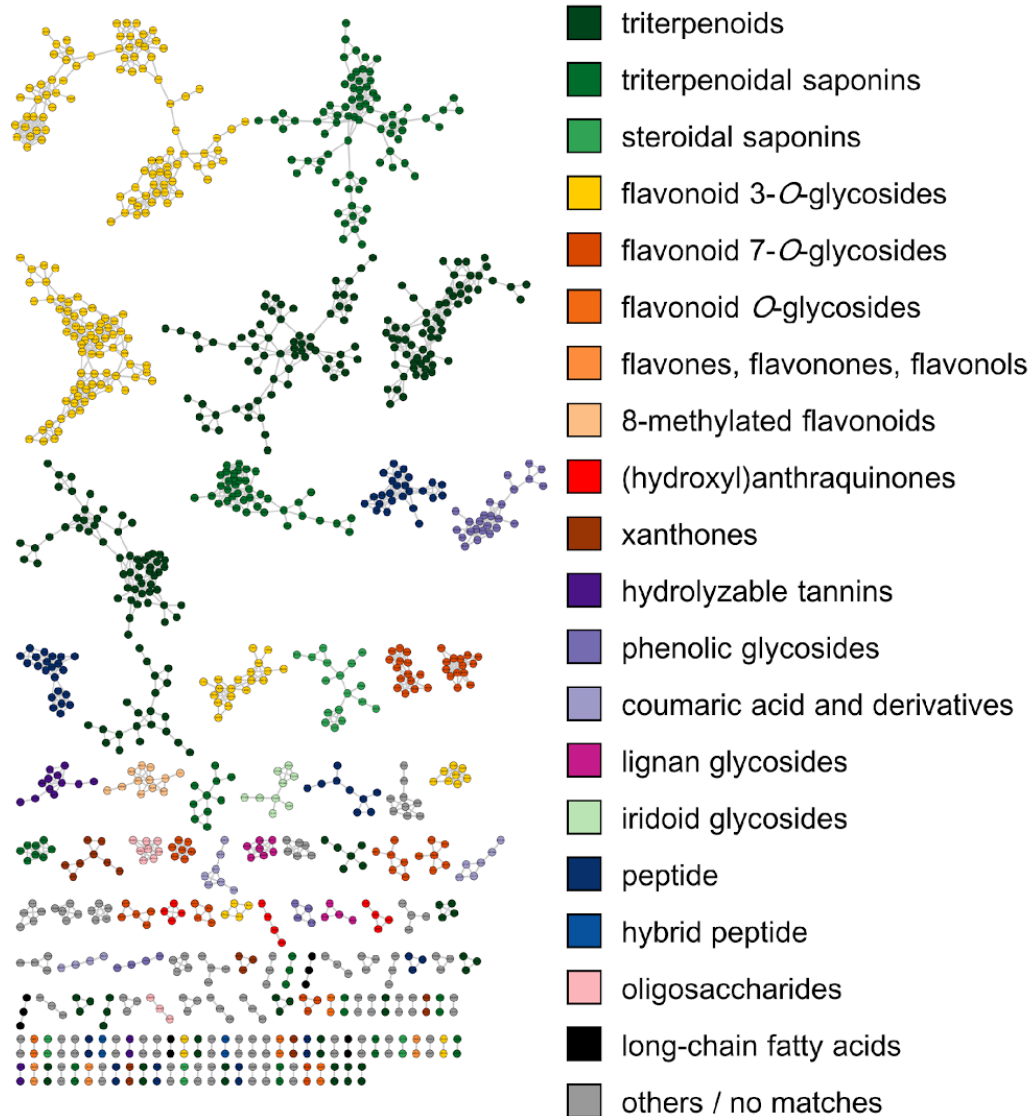
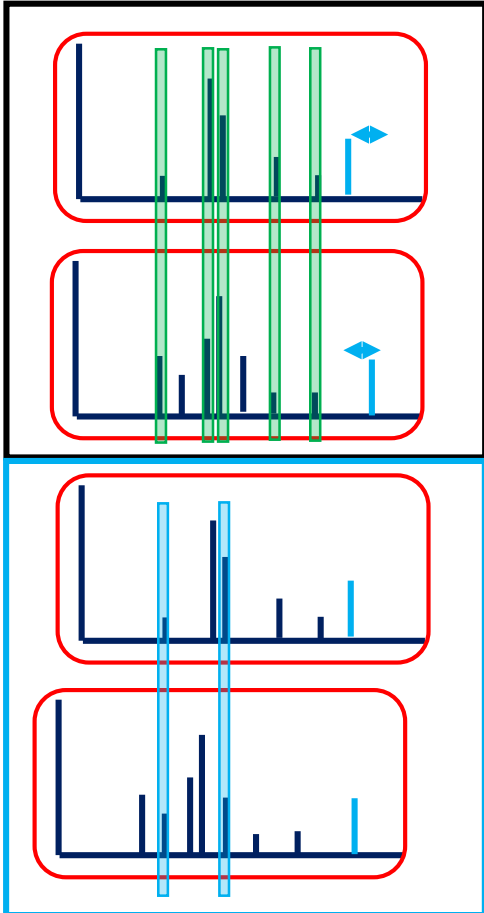
Dr Madeleine Ernst, UCSD





# Illuminating the Rhamnaceae plant chemistry

## Molecular Networking



## plant related classifications:

different flavonoids

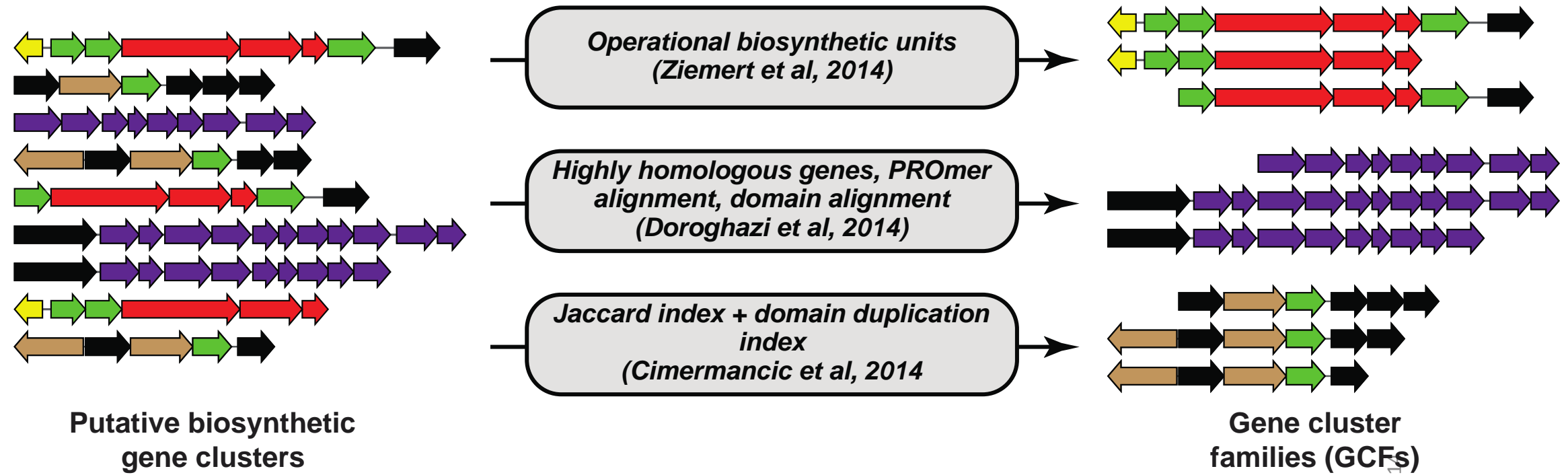
phenolic glycosides

triterpenoids

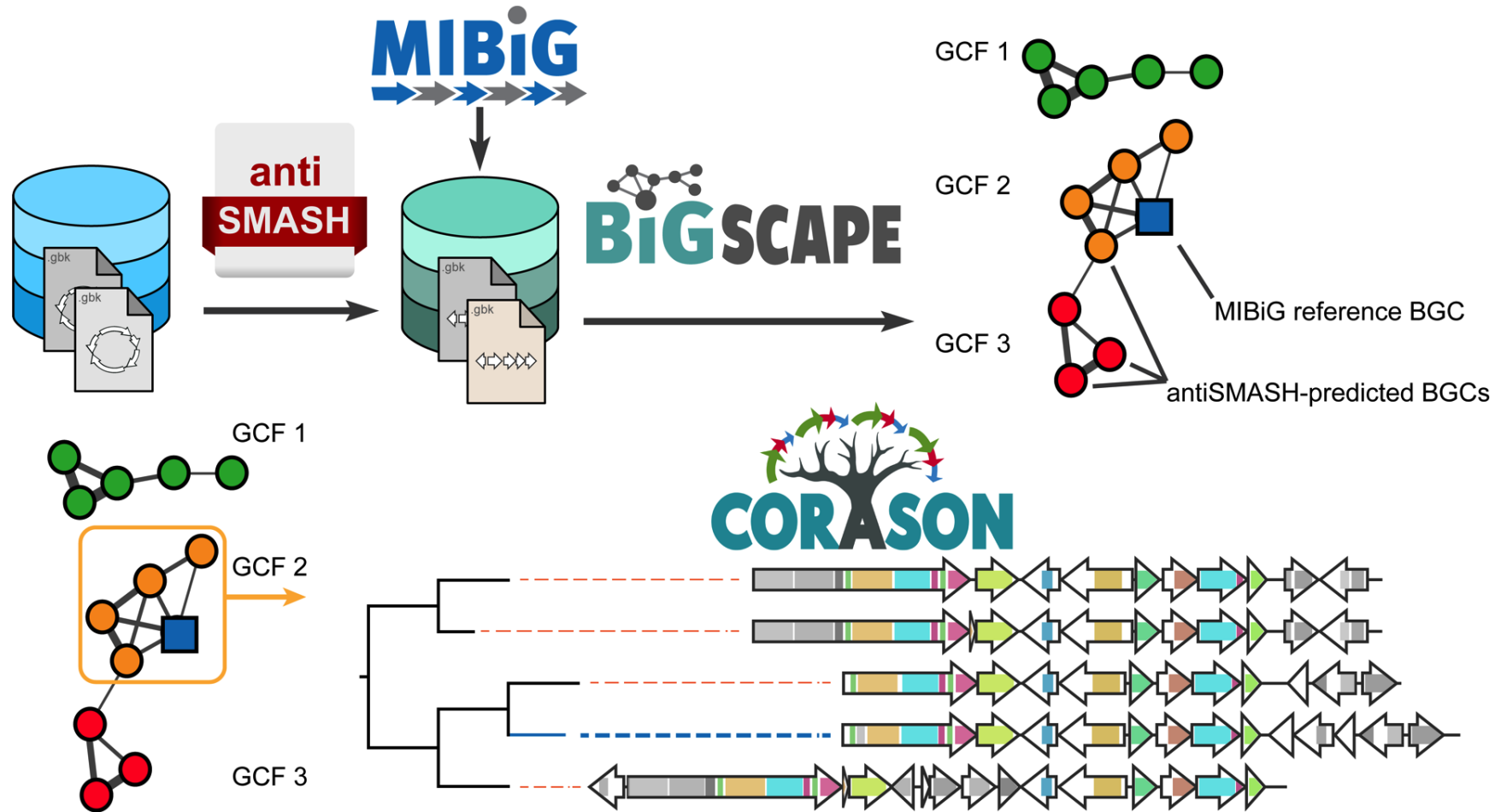
Dr Kyo Bin Kang, UCSD



# Grouping Biosynthetic Gene Clusters into Families



# Automated reconstruction and phylogenomic analysis



Dr. Jorge Navarro Muñoz



Nelly Selem-Mojica

# A FAIR Future Outlook

Lots of potential but no computer-readable formal linking options available

Community-driven platform:

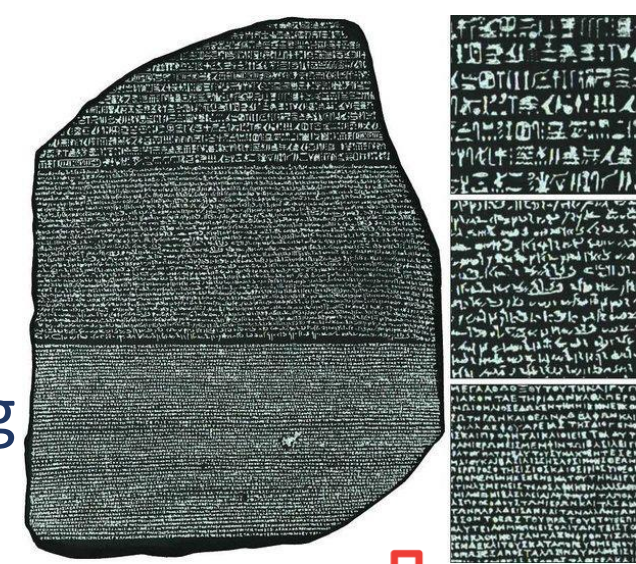
Collecting first feedback by close collaborators

**Minimum information needed versus time commitment**

The next step:

In collaboration with Glasgow we work on NPlinker:

**Systematic exploration of linking algorithms**



Paired Data Platform	
<u>Genome</u>	<u>MS/MS</u>
Name	MASSive /
NCBI Acc. #	GNPS link
Taxonomy ID	Experimental
antiSMASH	Parameters