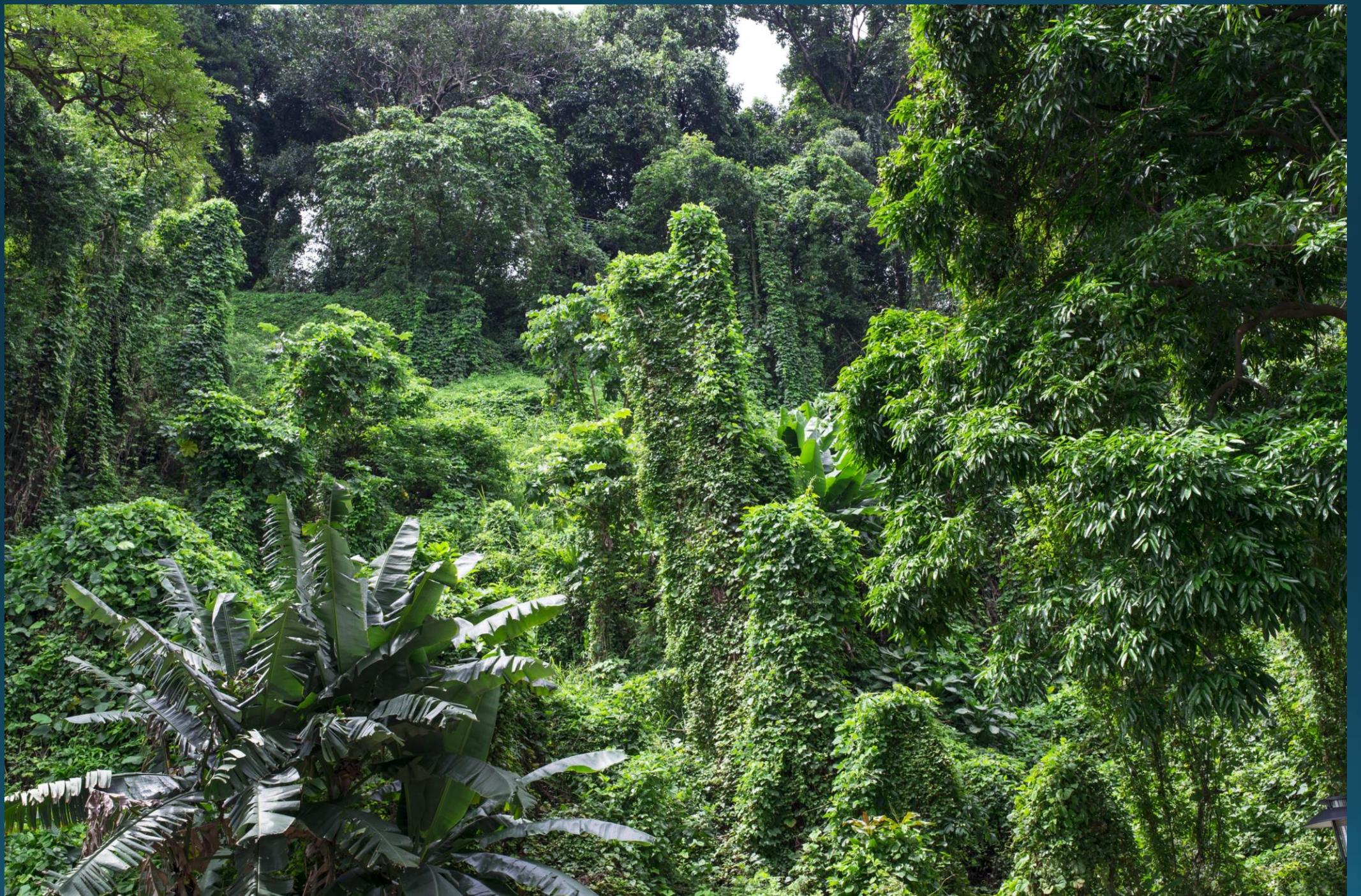


Cheaper by the Million:  
Expanding Untargeted Metabolomics  
into the  
Theoretical Realm

Anthony Mills

Coley/Kursar Laboratory

University of Utah Department of Biology



# Coley/Kursar Lab Question:



Why is the rainforest green?



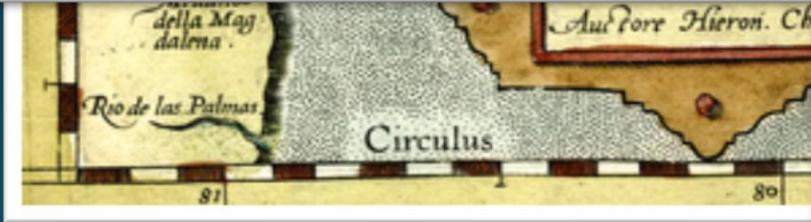
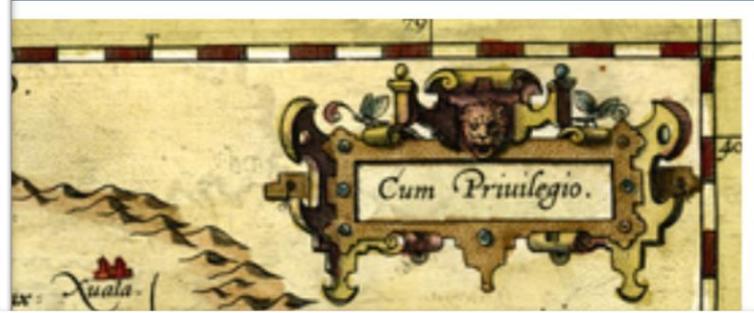








Juan Ponce de León



# *Hippomane mancinella*



# “La manzanilla de la muerte”





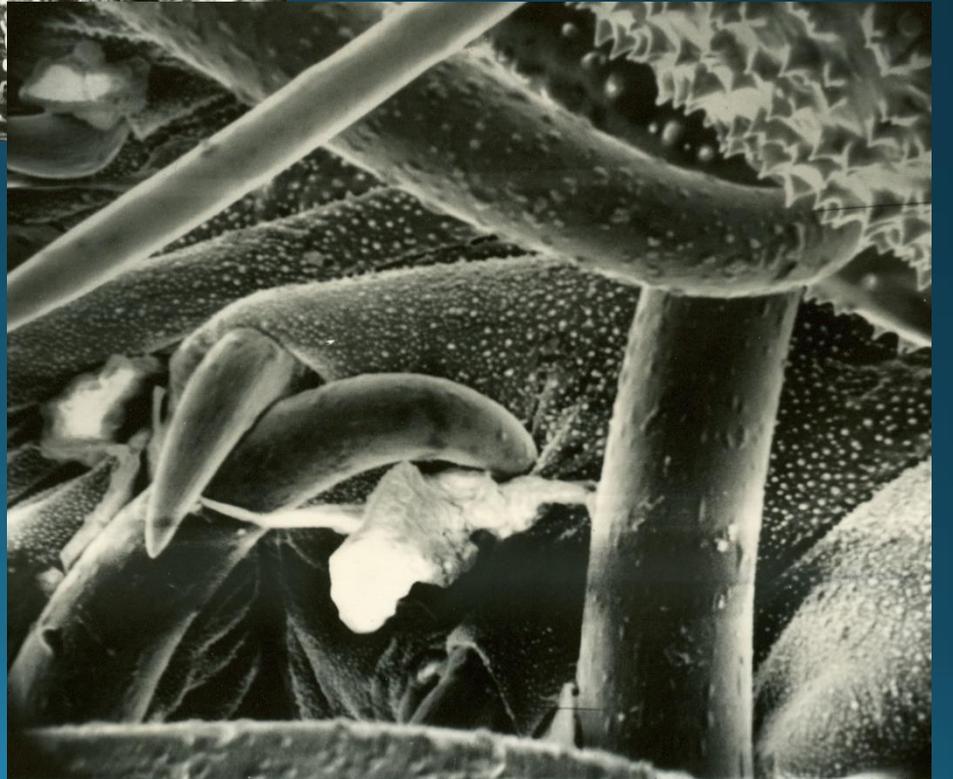
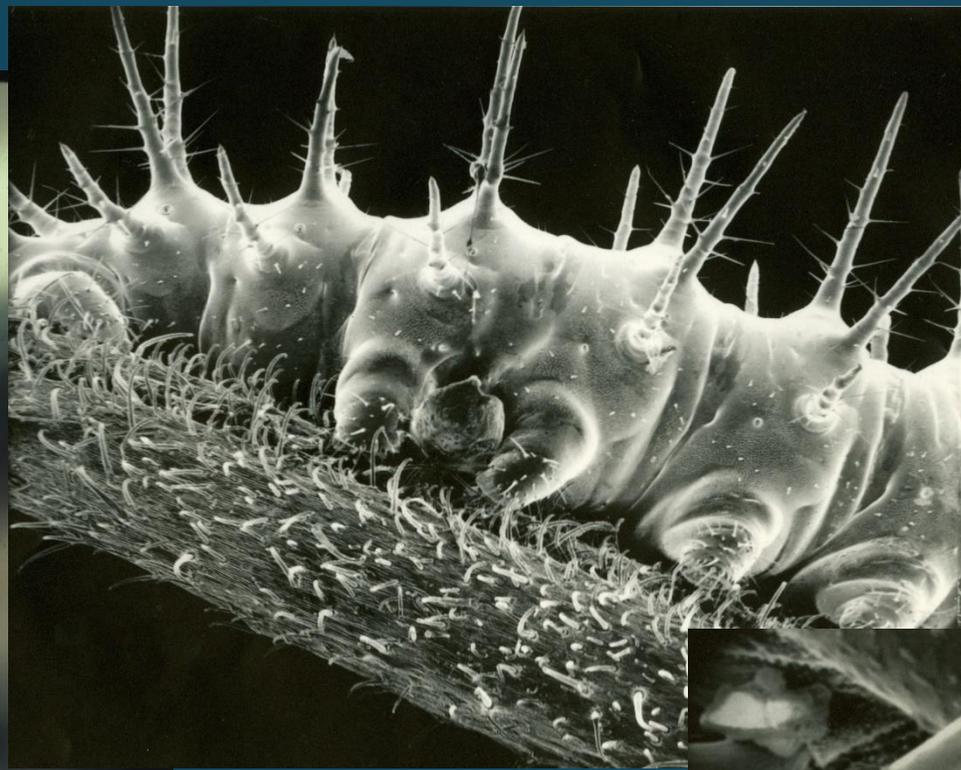


Who's the real enemy?





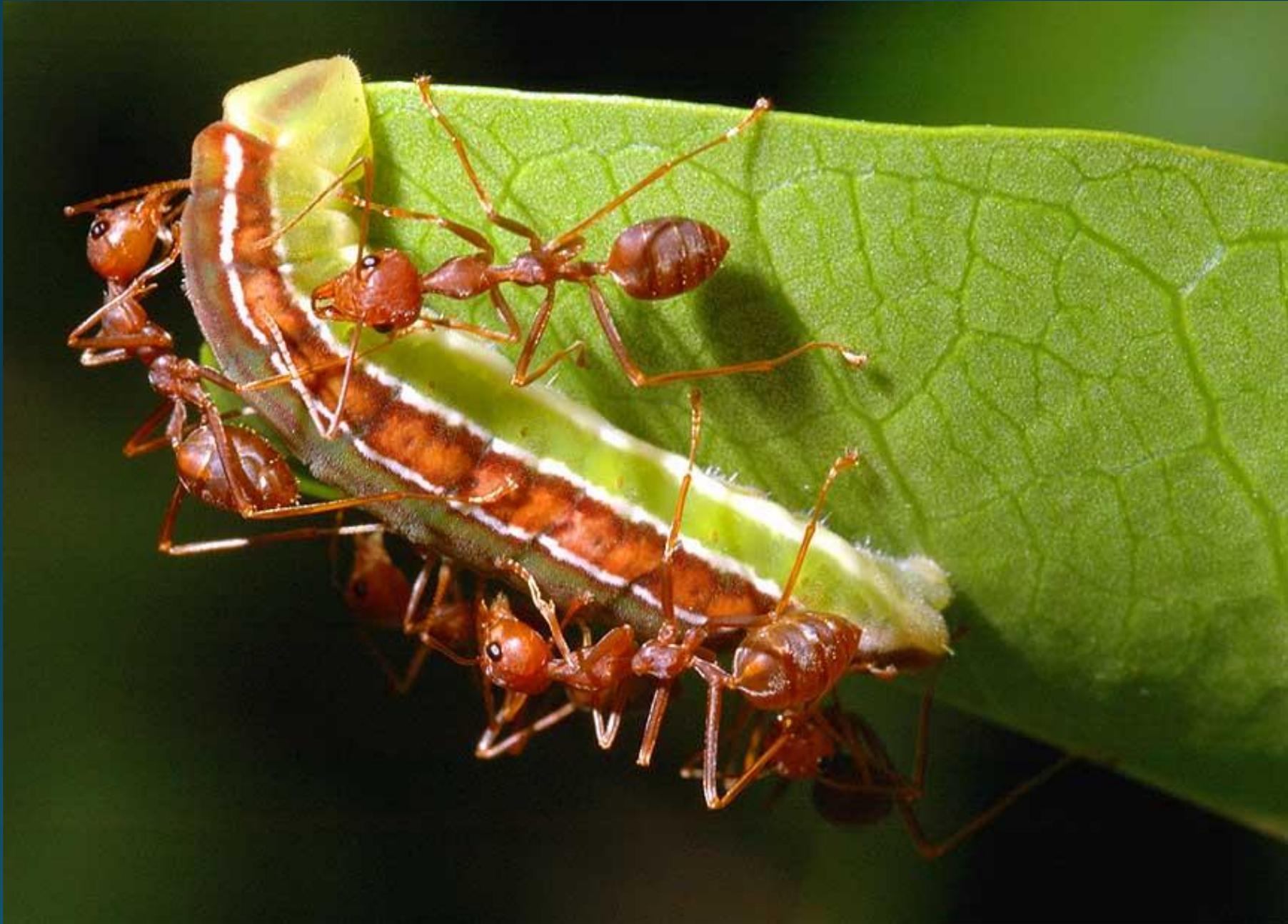


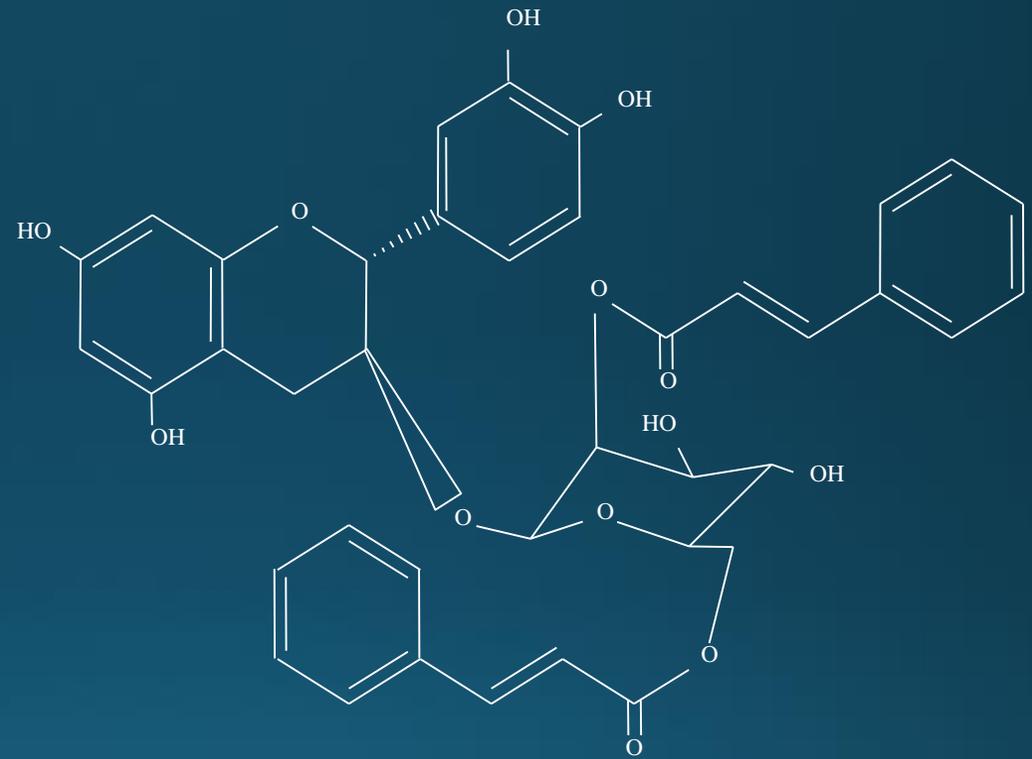




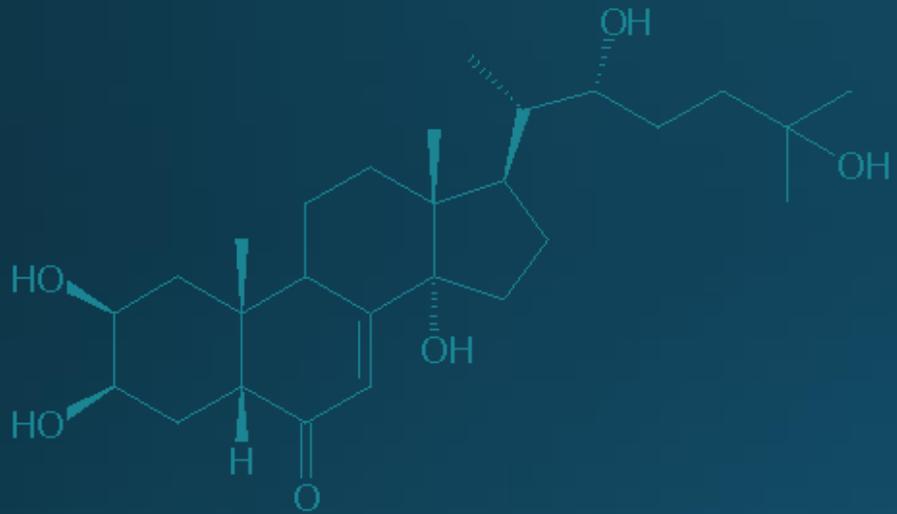
© Jeff Cremer / @JCremerPhoto

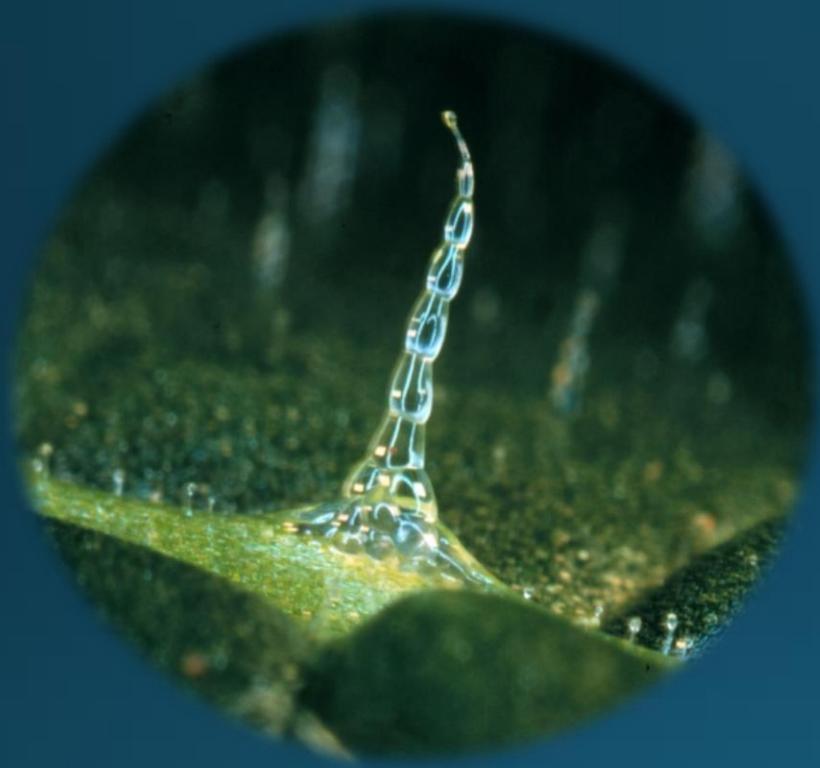






# Ecdysone









# Plant-Herbivore Interactions

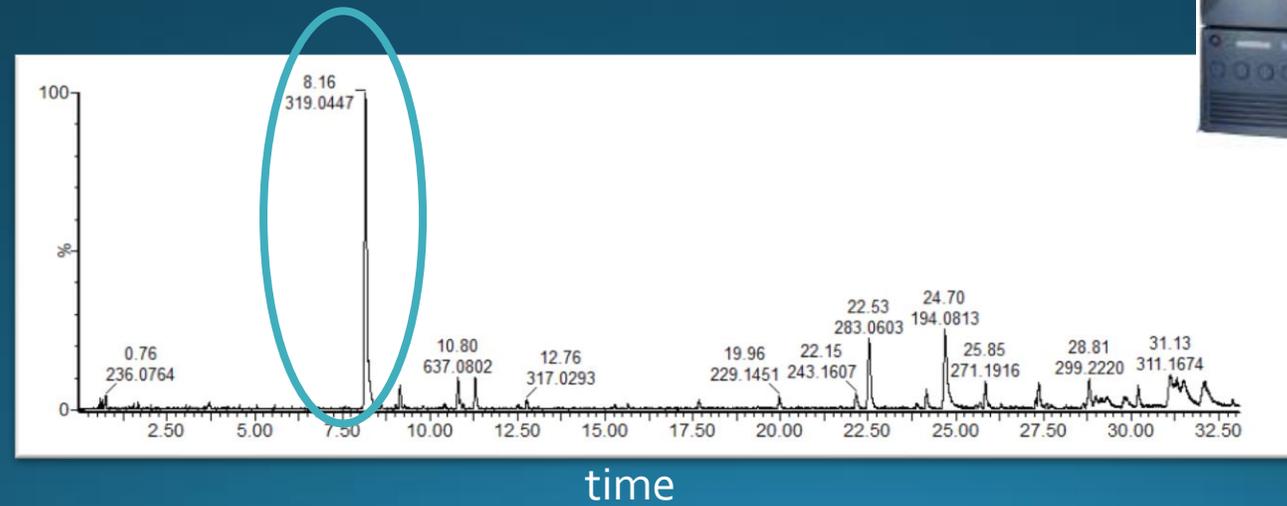
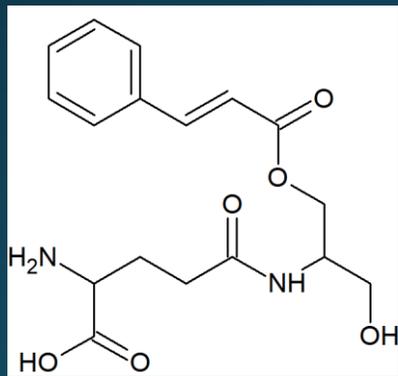
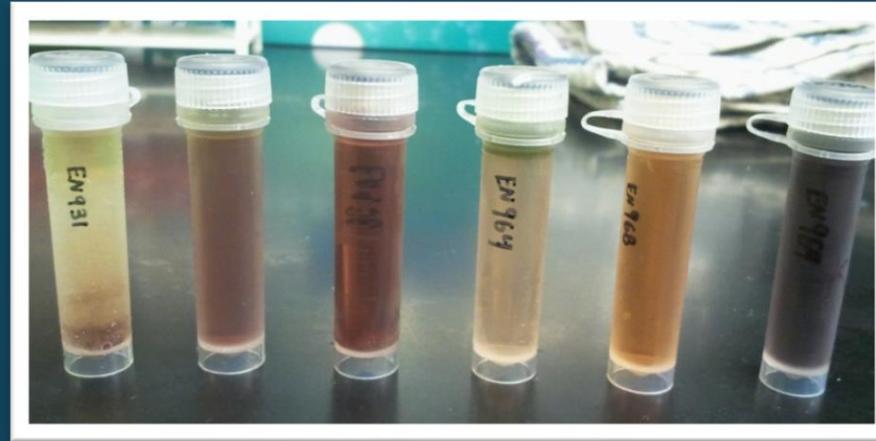




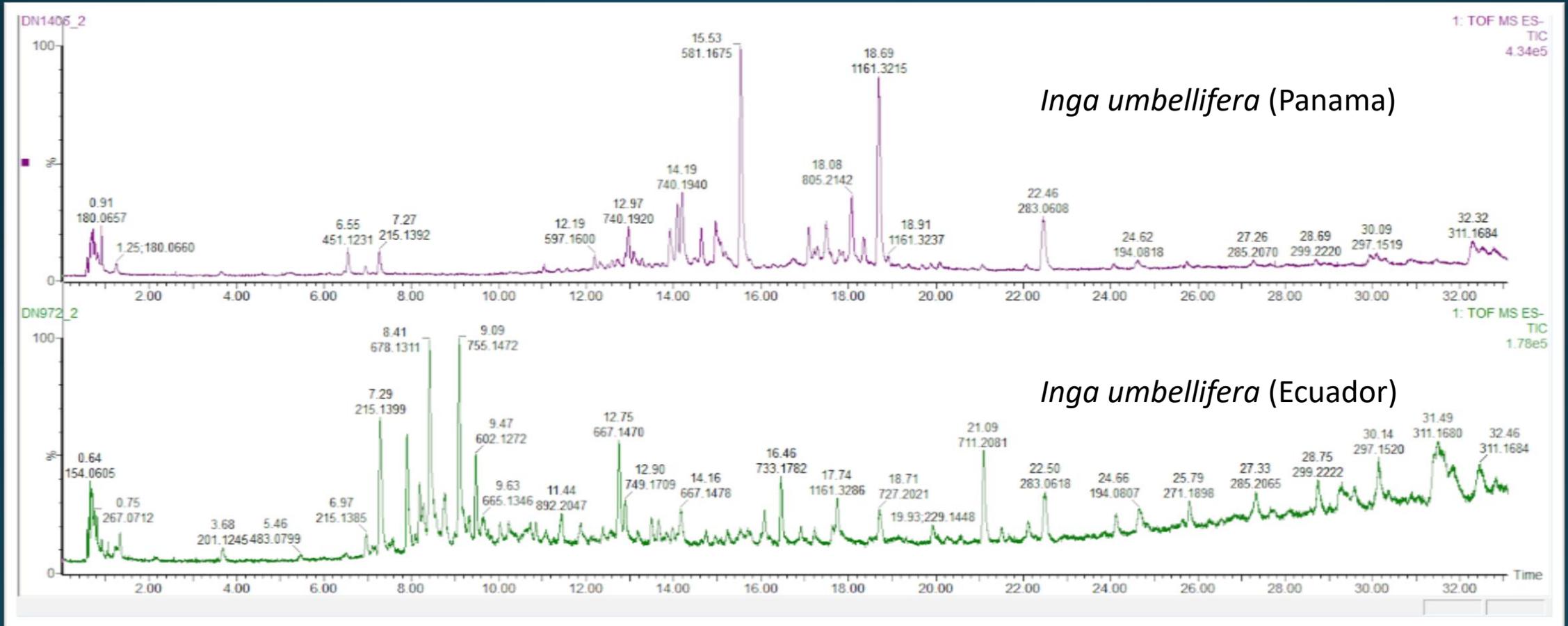
# *Inga*

- Neotropical genus
- Legume (Fabaceae)
- What's in them?
  - 3,600 Compounds
    - Only 4% known
  - Ecologically important
  - Central goal of lab

# Procedure



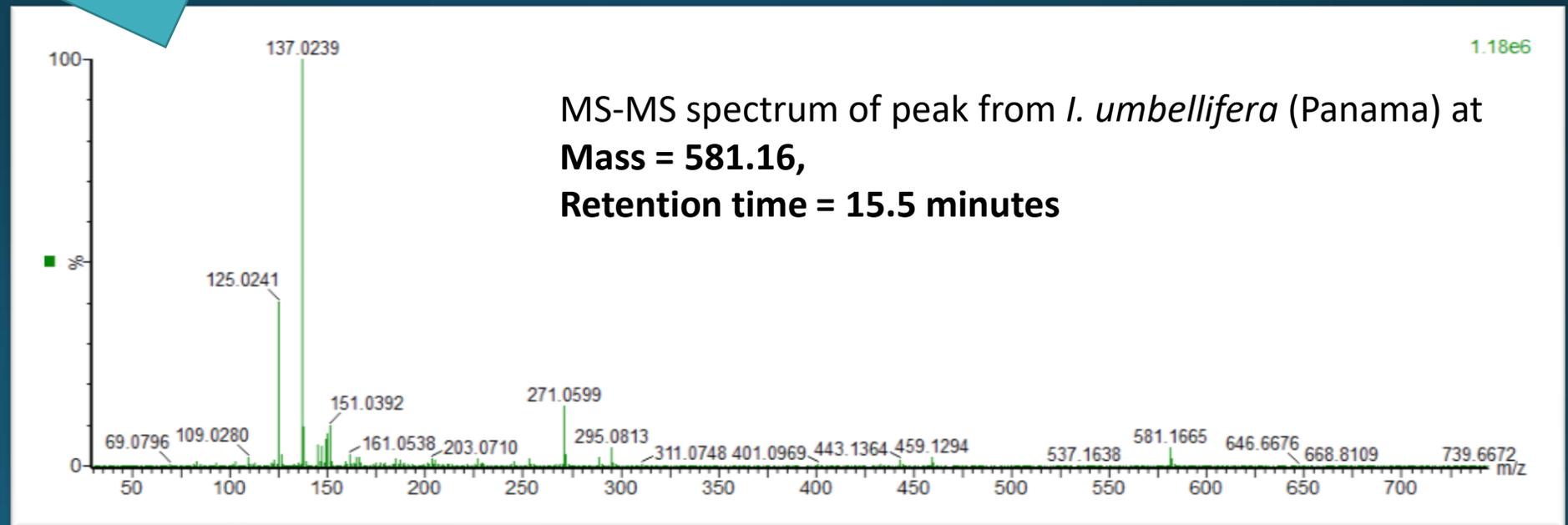
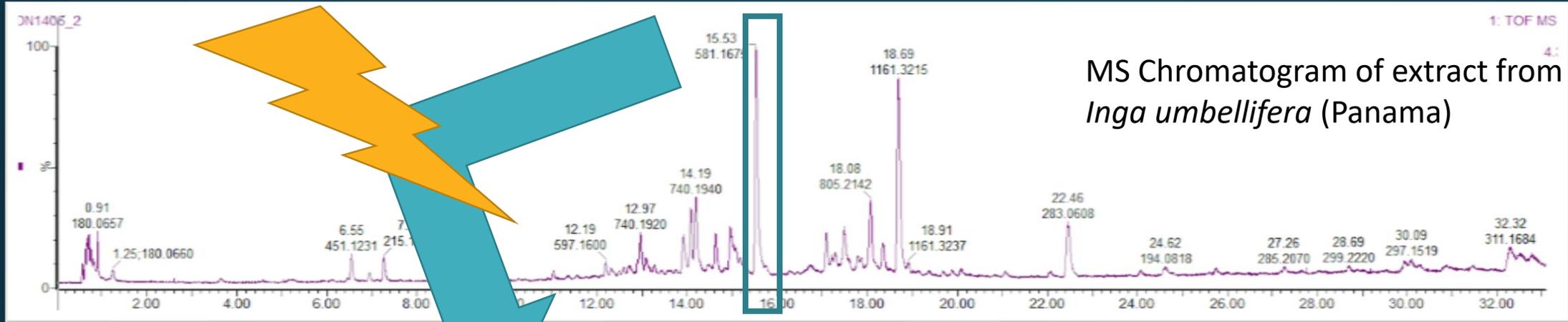
# Chemotypes



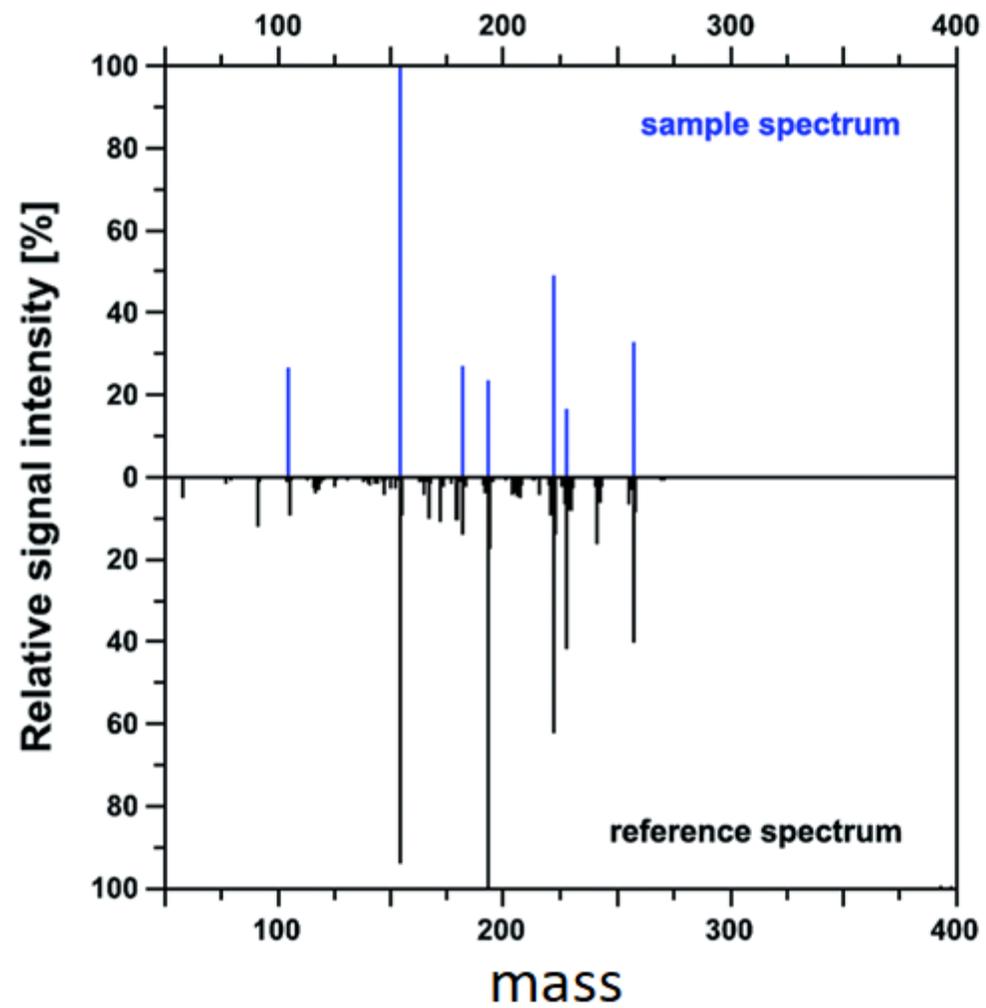
# Goal: Improve Identification

- Purification and NMR are costly and time+labor intensive
- Untargeted metabolomics still a developing field
- High-Res Mass is starting point
- MS-MS gives structural information

# Tandem Mass Spectrometry



# Diazepam



- NP Databases do exist
- Few if any MS-MS spectra
- Structures present

# *In-Silico* fragmentation



☆ Utilities ▾

Help

Data

Publications

Contact Us

## CFM-ID

Competitive Fragmentation Modeling for Metabolite Identification

### Spectra Prediction

Predicts the spectra for a given input molecule. Spectra are computed for low (10V), medium (20V) and high (40V) collision energy levels and are represented by a list of 'mass intensity' pairs, each corresponding to a peak in the spectra.

**Parent Compound Structure**

InChI or SMILES format

Enter an InChI or SMILES string

InChI strings need to start with "InChI=" and are not expected to have any charge - an additional H+ will be added. Maximum compound size is 200 atoms. Load an [InChI example](#), [SMILES example](#), or another [SMILES example](#).

**Spectra Type**

ESI ▾

**Ion Mode**

Positive ▾

**Adduct Type**

[M+H]<sup>+</sup> ▾

Reset

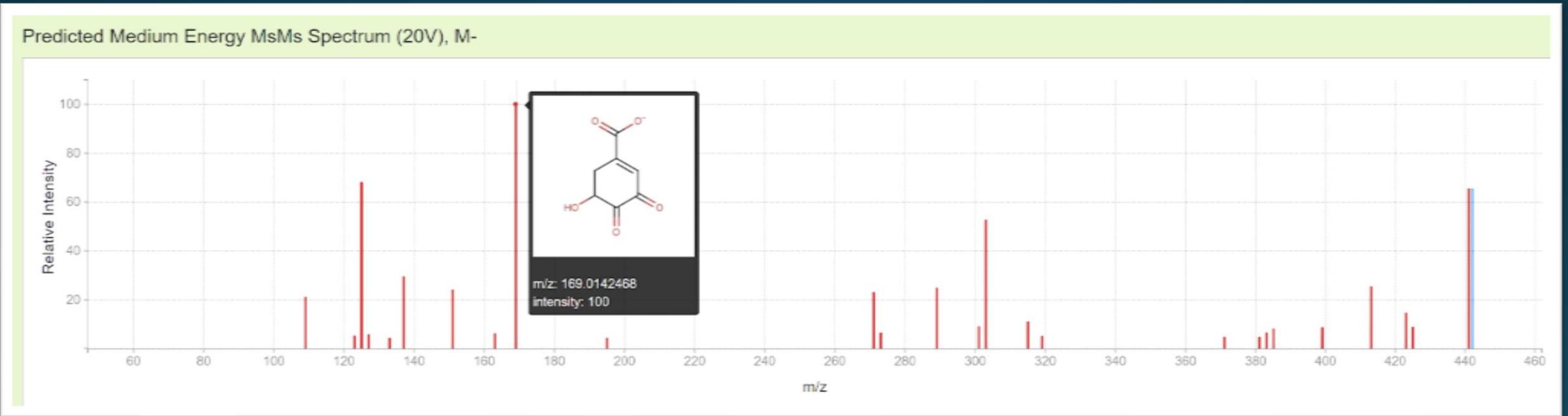
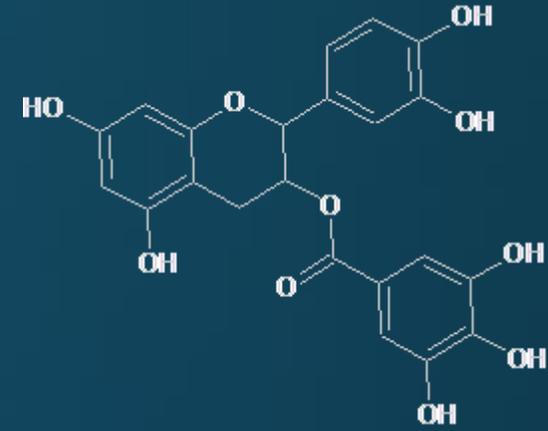
Submit

If you wish to run multiple jobs, input larger query molecules, or customize the computation parameters, you can freely download the source code here: <http://sourceforge.net/projects/cfm-id>.

# *In-Silico* fragmentation

## Epicatechin Gallate:

O=C(O[C@@H]2Cc3c(O[C@@H]2c1ccc(O)c(O)c1)cc(O)cc3O)c4cc(O)c(O)c(O)c4



~230,000 Compounds



# CHPC

## Ember Cluster

- 144 Dual Socket-Six Core Nodes (1728 total cores)
- 2.8 GHz Intel Xeon (Westmere X5660) processors
- RAM:
  - 24 Gbytes memory per node (2 Gbytes per processor core) on most nodes





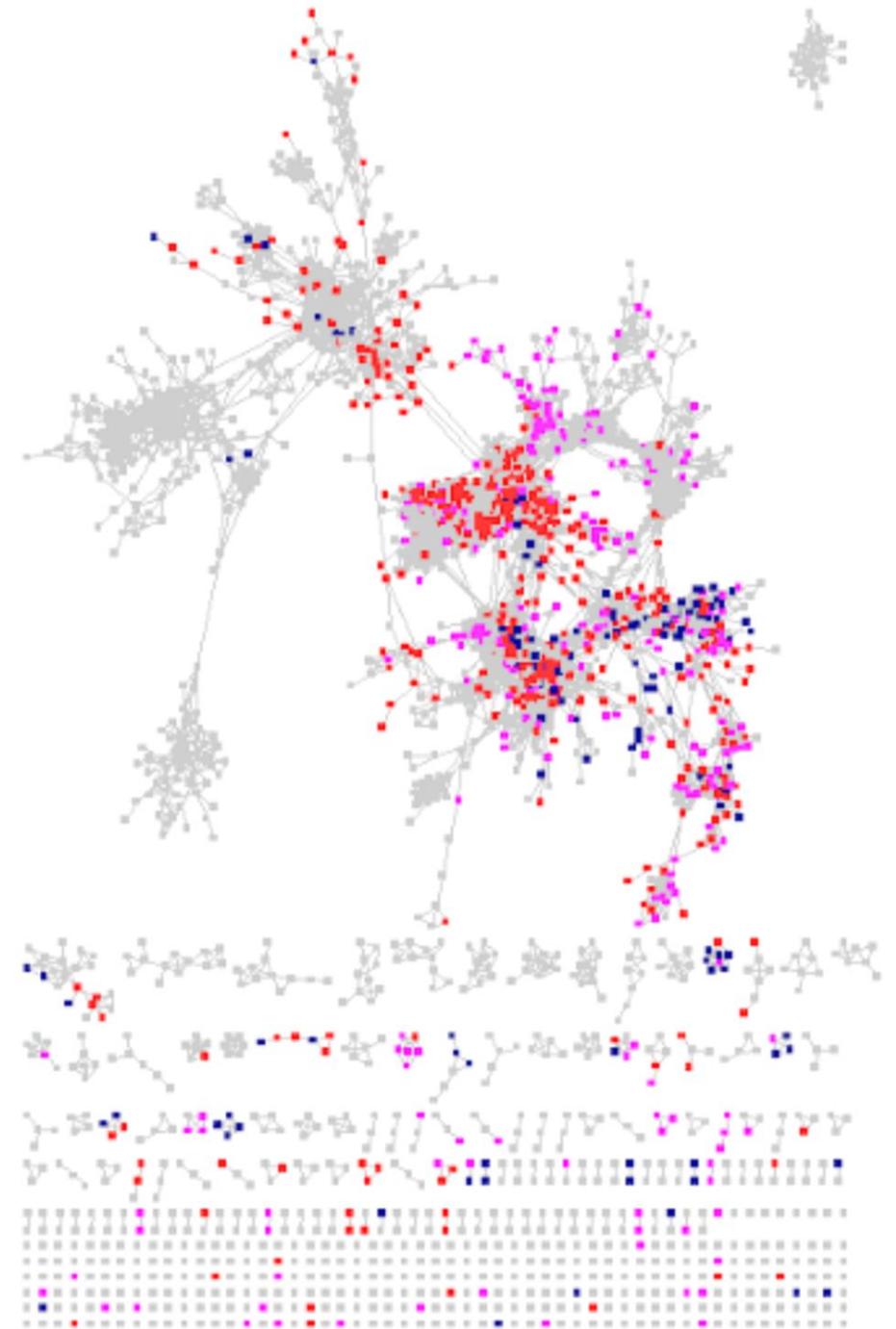
**Open Science Grid**

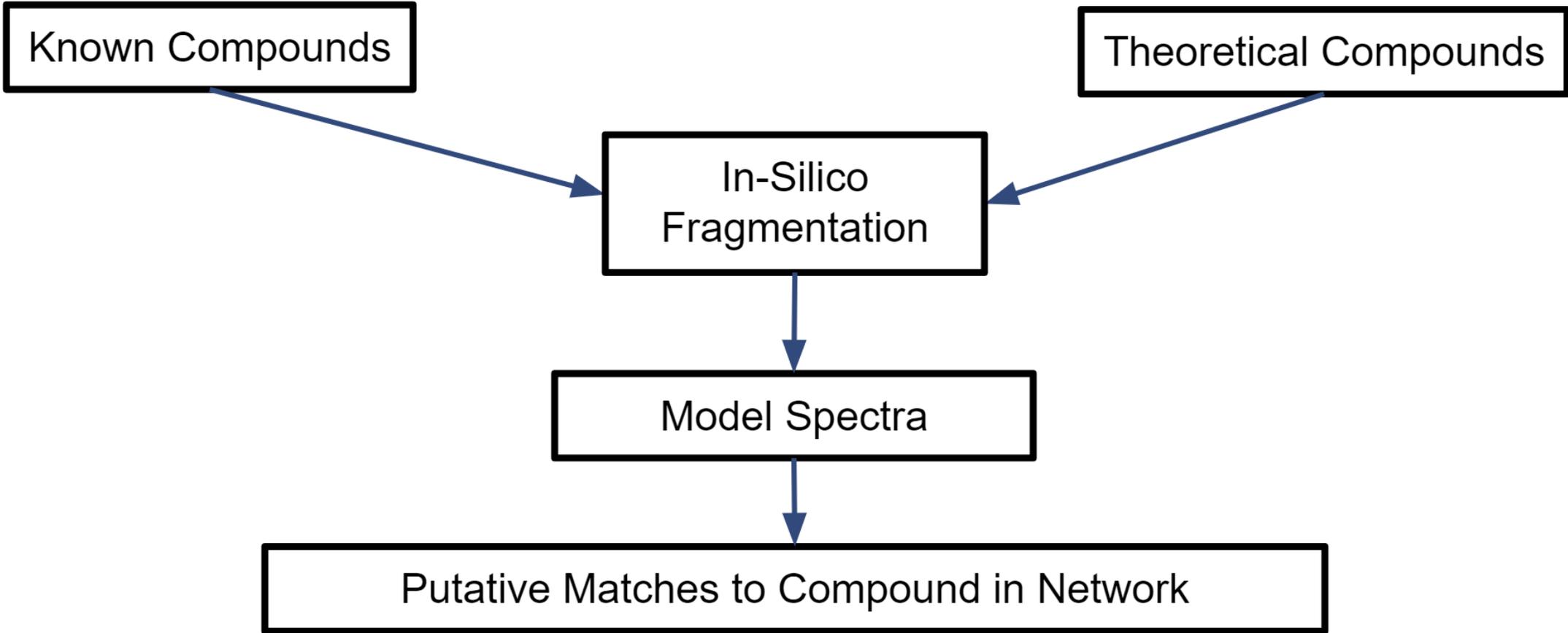
# Comparison of HTC to HPC

<b>230,737 molecules total</b>	<b>Job Size</b>	<b>Run Rate</b>	<b>Free Nodes</b>	<b>Time to Completion</b>
<b>Ember Cluster CHPC</b>	462 files containing 500 molecules each	69 molecules / hour	10-20	334 hours (based on 10 free nodes)
<b>OSG</b>	4,614 files containing 50 molecules each	4 molecules / hour	Thousands	12 hours

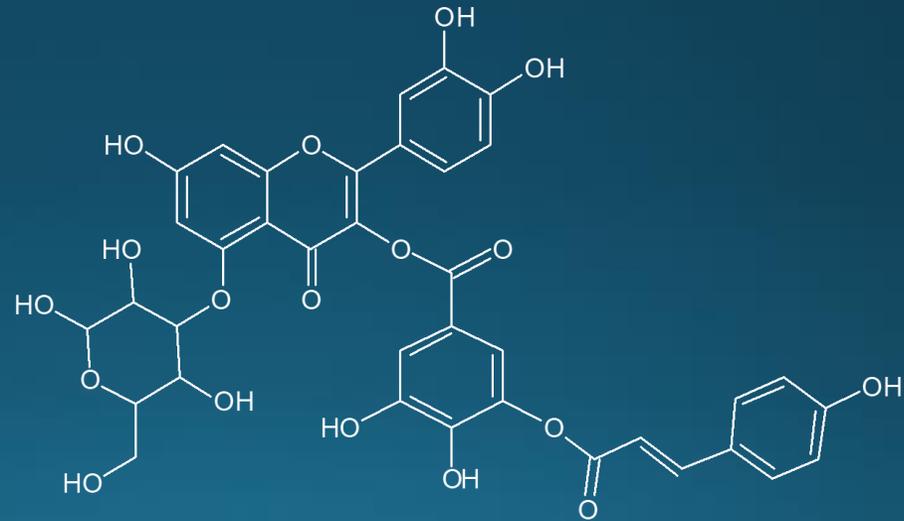
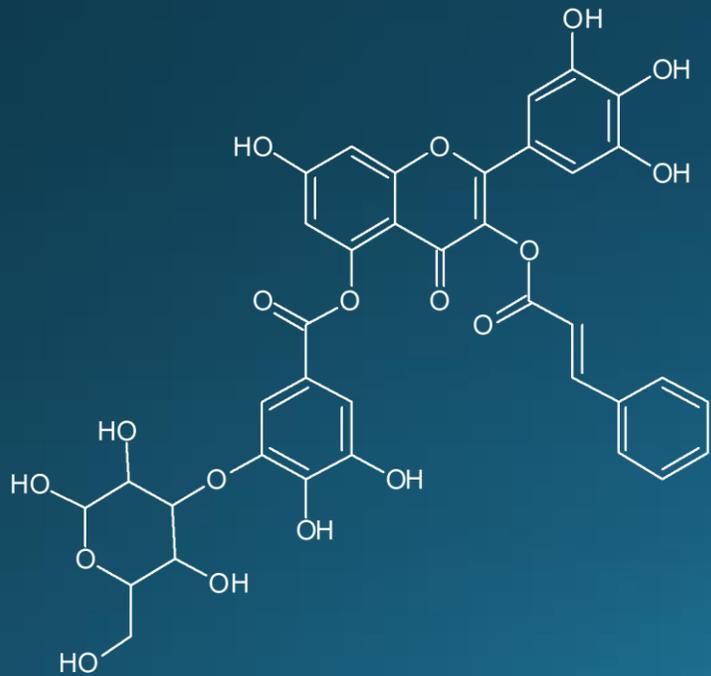
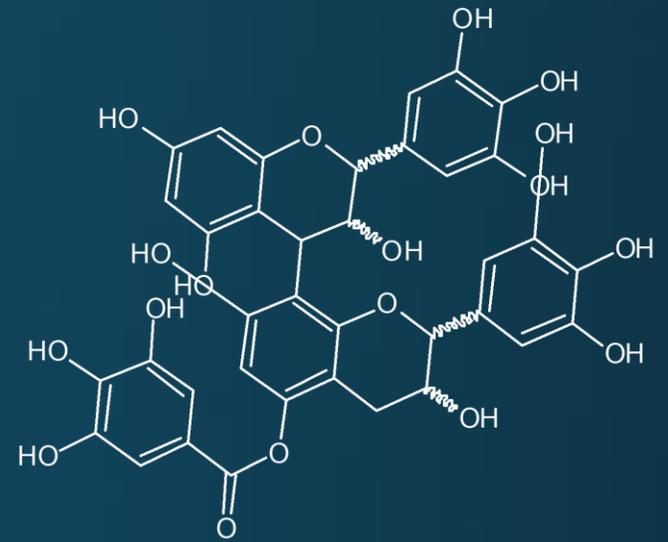
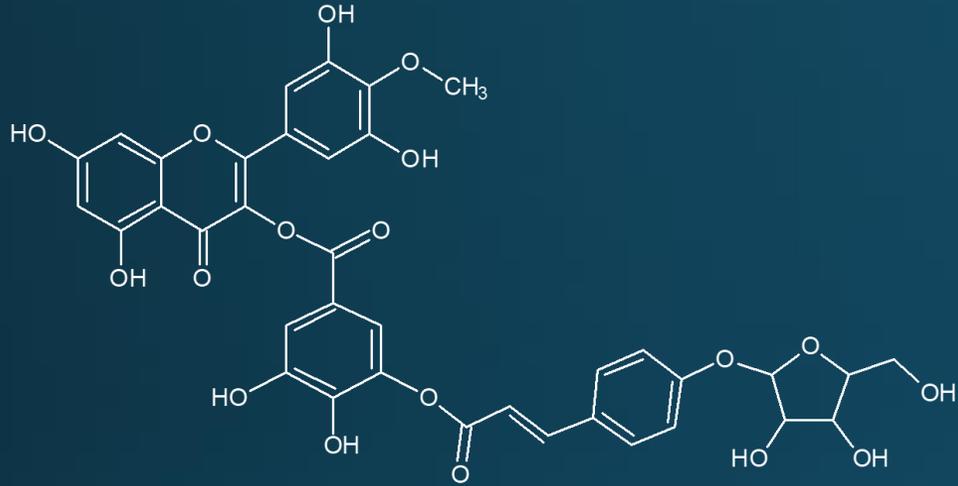
# Results:

- **Blue** = match to Library MS-MS spectrum
- **Red** = match to *in-silico* MS-MS spectrum
- **Magenta** = match to mass only

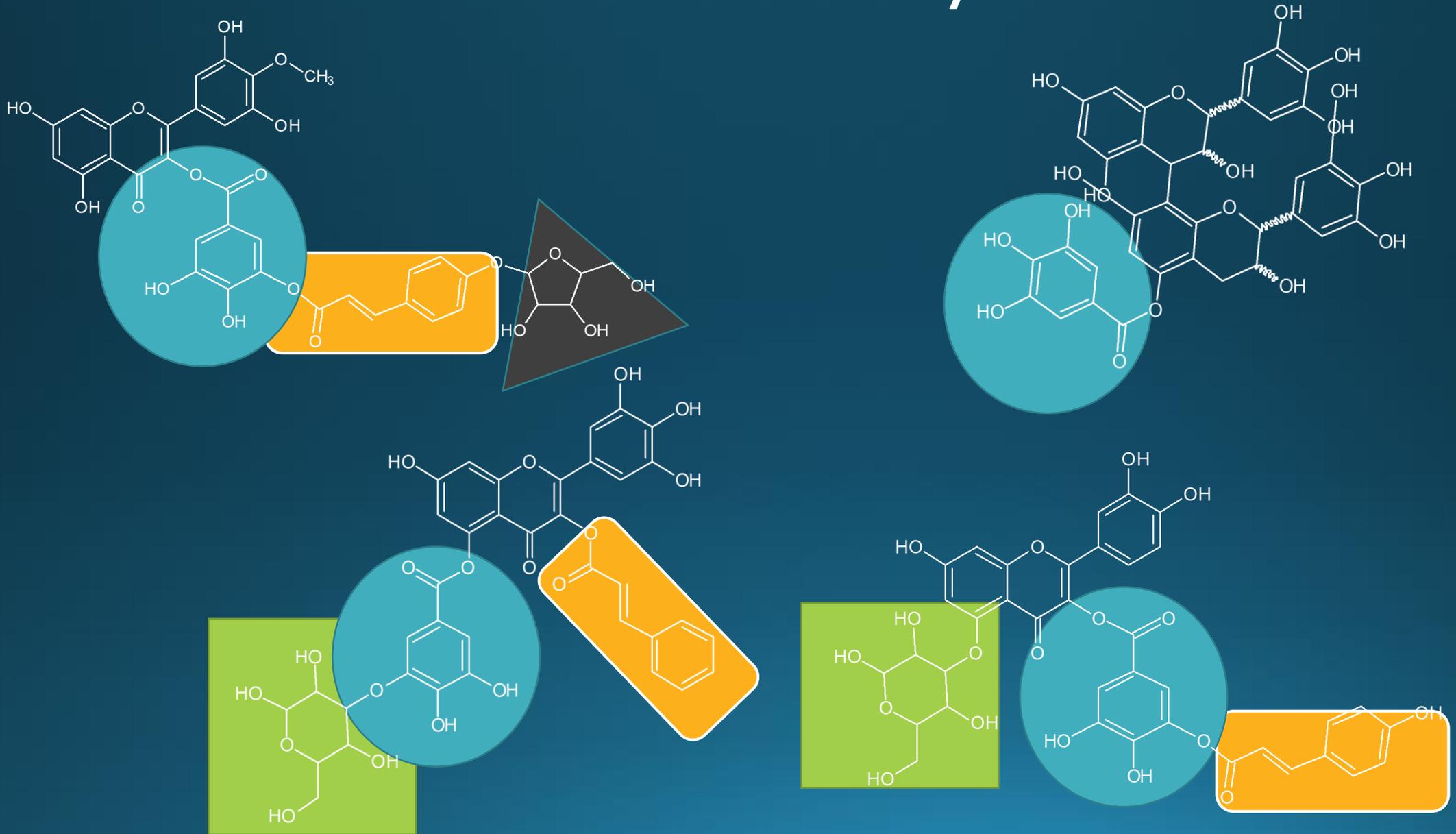




# Combinatorial Chemistry



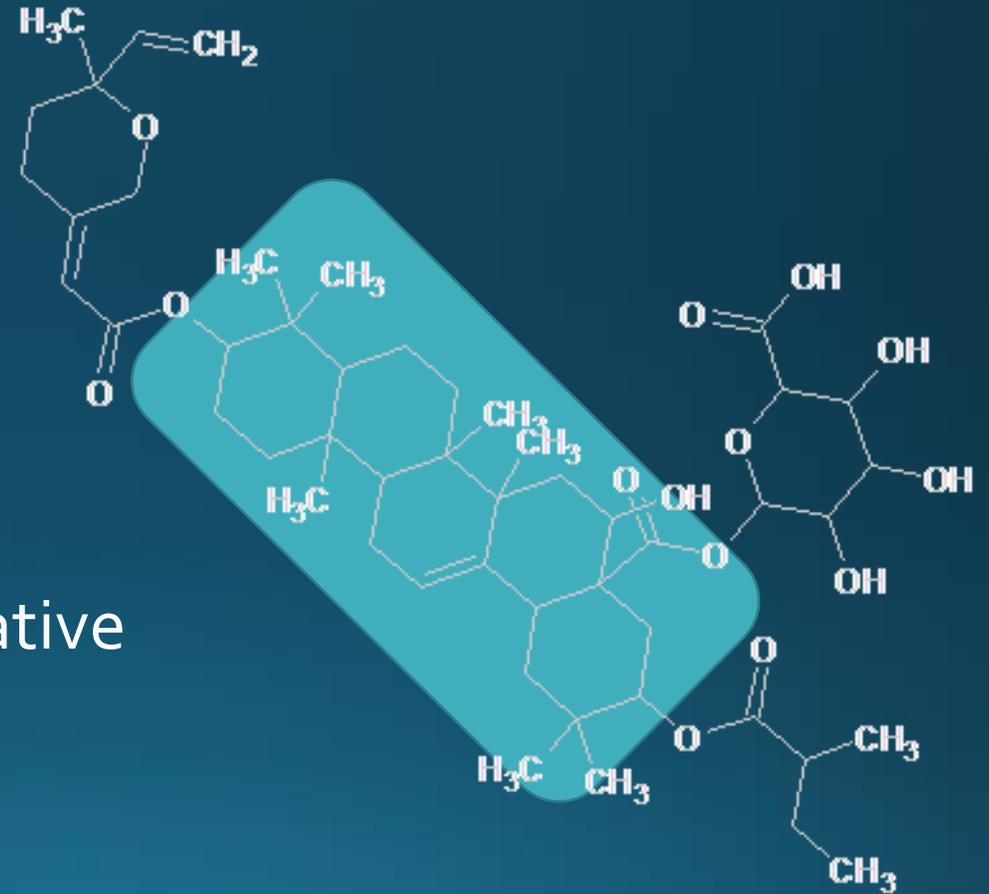
# Combinatorial Chemistry





# Theoretical Saponins

- 7 Saponin Cores
- 21 Substituents
- Fully enumerated library
  - **4.2 million compounds...**
  - Reduced to 80,000
    - Fingerprints are not representative
  - Scaling back up



# Acknowledgements

- Coley/Kursar Lab

- Thomas Kursar
- Lissy Coley
- Dale Forrister
- Gordon Younkin

- University of Utah CHPC

- Wim Cardoen
- Anita Orendt
- Martin Čuma

- Open Science Grid

- Bala Desinghu
- Tim Cartwright
- User School Personnel



Open Science Grid



PEKING  
UNIVERSITY



Department of Biology

Bioscience Undergraduate Research Program

Center for High Performance Computing