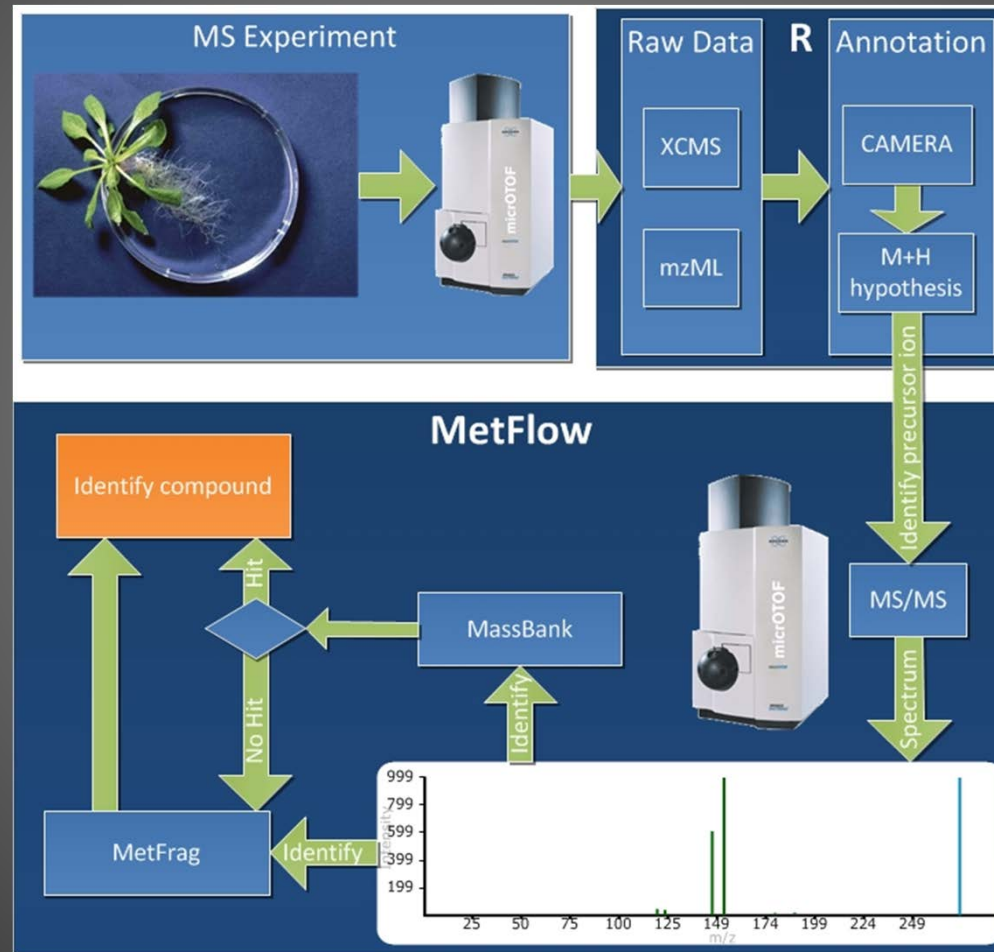# MetFusion: integration of compound identification strategies

## MassBank Workshop 27.11.2012

## Amsterdam, The Netherlands

Michael Gerlich, Department of Stress and
Developmental Biology, IPB Halle

# Introduction

Michael Gerlich, Department of Stress and
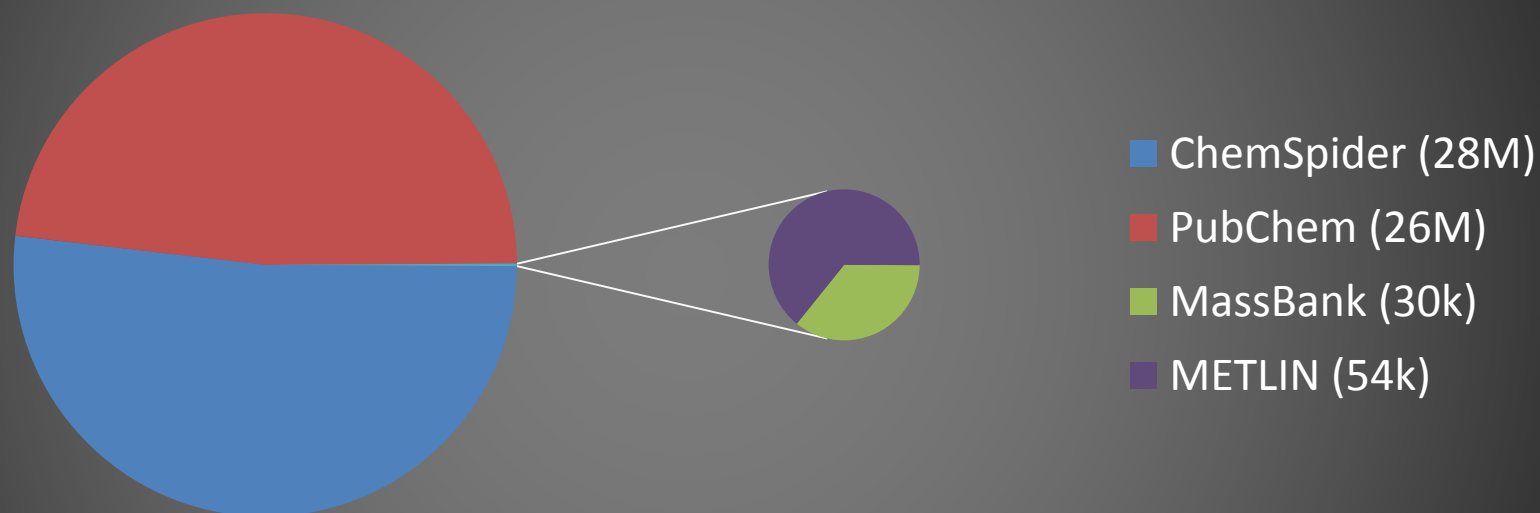Developmental Biology, IPB Halle

# Structures & Spectra

- Target list of interesting MS$^2$ spectra
- Requires expert knowledge
  - Time-consuming
  - Impossible to keep track of high-throughput
- Only small fraction of compounds has associated reference spectrum
  - Required for speed-up in re-identification
  - Difficult to use for *de novo* identification

# Structures vs. Spectra Imbalance

**Database Entries**



- ChemSpider (28M)
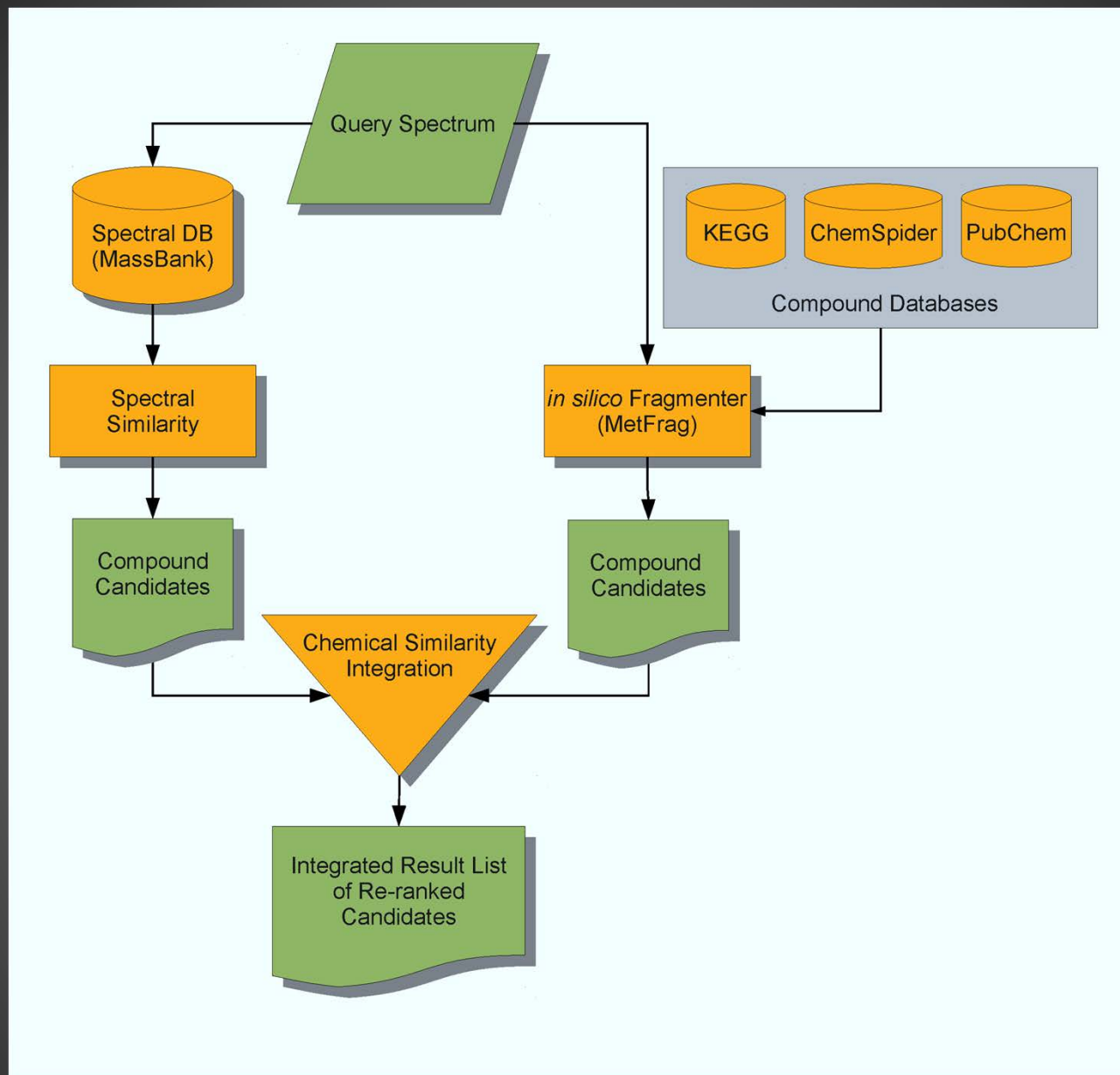- PubChem (26M)
- MassBank (30k)
- METLIN (54k)

- Millions vs. thousands, overlap unknown
- Spectrum queries not possible for compound databases

# Observation vs. Prediction

- Observation: MassBank
  - Search measured spectra with peak list
  - Find compounds with matching/similar spectrum
  - Few reliable results

- Prediction: MetFrag
  - Combinatorial Fragmenter
  - Generate fragments *in silico*, matches to peak list
  - Uses compound databases,
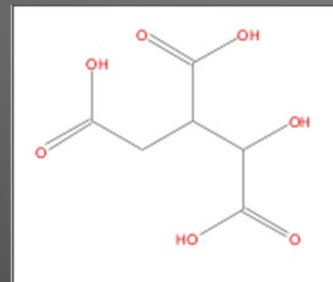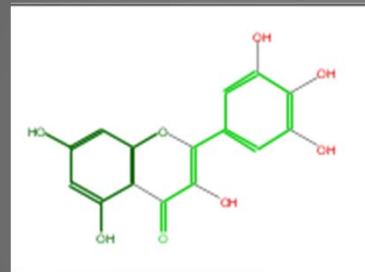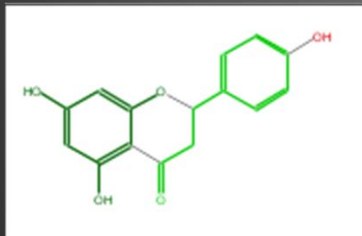    many possible predicted results

# MetFusion = Observation, Prediction & Similarity

- Combine results via chemical similarity
  - Structural fingerprints
  - Use each information (score)
  - Avoid strict limits/thresholds
- Aim: improve identification
  - Assume that correct compound is present in compound database (larger coverage)
  - Enhance MetFrag results with spectral data
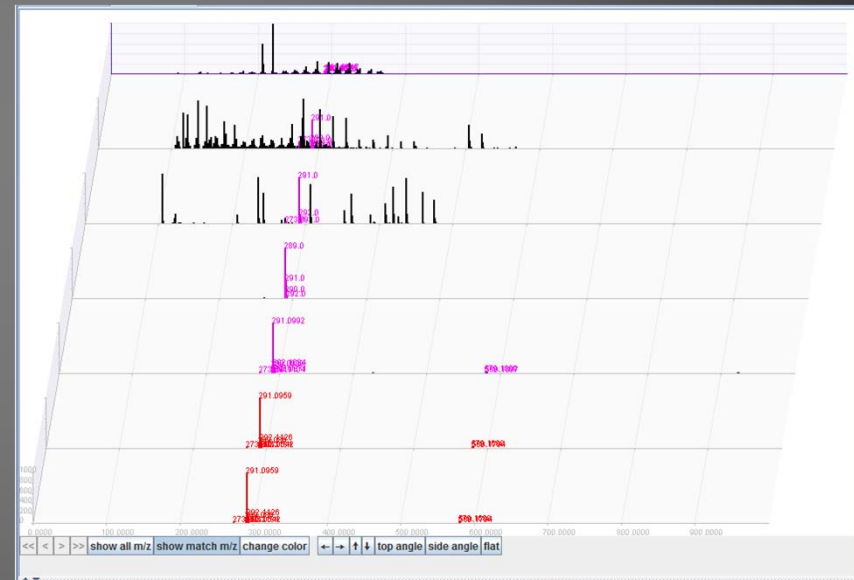  - DOI 10.1002/jms.3123

# What is similar?

**Chemical Similarity**

**Spectral Similarity**

# Similarity Measures

## Chemical Similarity

- Tanimoto coefficient matches properties

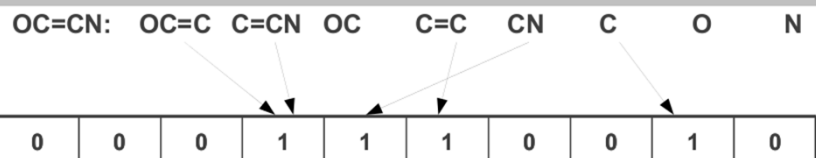- $Tan = \frac{C}{A+B-C}$

- [0,1]

## Spectral Similarity

- Modified Cosine distance

- Matches peak masses & intensities

- $W_i = int_i^m * m/z_i^n$

- [0,1]

Path based (FP2): Paths of up to *n* atoms are generated and hashed to set the bits

OC=CN:   OC=C   C=CN   OC      C=C     CN      C       O       N

| 0 | 0 | 0 | 1 | 1 | 1 | 0 | 0 | 1 | 0 |

Dr. Ernst-Georg Schmid Universität Duisburg
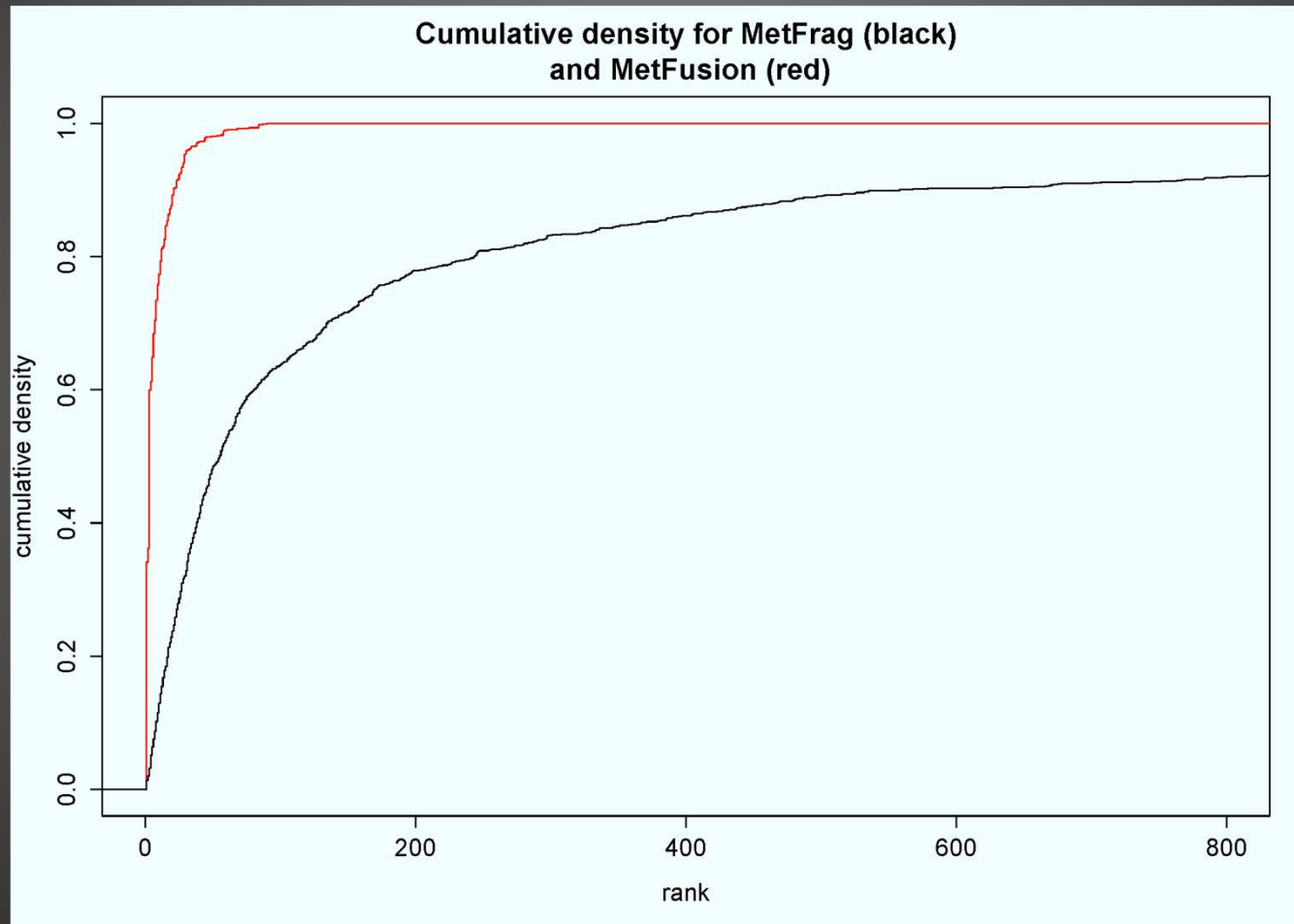GCC 2010

# Similarity Matrix

| MF\MB | C00509[0.975] | C06561[0.965] | C09099[0.956] | C09789[0.916] | C03406[0.599] | C04577[0.520] | C00158[0.502] | C10107[0.468] | C00311[0.418] | -----[0.413] |
|---|---|---|---|---|---|---|---|---|---|---|
| C00509[1.000] | 1 | 0,299 | 0,721 | 0,632 | 0,14 | 0,152 | 0,106 | 0,464 | 0,106 | 0,338 |
| C16232[1.000] | 0,916 | 0,293 | 0,687 | 0,617 | 0,14 | 0,152 | 0,1 | 0,468 | 0,1 | 0,365 |
| C06561[0.966] | 0,299 | 1 | 0,252 | 0,243 | 0,102 | 0,142 | 0,097 | 0,445 | 0,097 | 0,259 |
| C12087[0.966] | 0,25 | 0,316 | 0,24 | 0,243 | 0,122 | 0,212 | 0,089 | 0,328 | 0,089 | 0,32 |
| C14458[0.966] | 0,618 | 0,316 | 0,5 | 0,45 | 0,113 | 0,149 | 0,091 | 0,38 | 0,091 | 0,289 |
| C09826[0.909] | 0,9 | 0,289 | 0,701 | 0,629 | 0,126 | 0,153 | 0,102 | 0,494 | 0,102 | 0,35 |
| C03567[0.462] | 0,582 | 0,316 | 0,479 | 0,442 | 0,11 | 0,149 | 0,088 | 0,379 | 0,088 | 0,292 |
| C09614[0.462] | 0,913 | 0,292 | 0,699 | 0,624 | 0,14 | 0,155 | 0,1 | 0,478 | 0,1 | 0,36 |
| C09751[0.443] | 0,904 | 0,292 | 0,704 | 0,632 | 0,132 | 0,152 | 0,102 | 0,504 | 0,102 | 0,354 |
| C09047[0.426] | 0,376 | 0,411 | 0,332 | 0,322 | 0,119 | 0,141 | 0,077 | 0,6 | 0,077 | 0,248 |
| C17673[0.426] | 0,355 | 0,323 | 0,322 | 0,303 | 0,133 | 0,12 | 0,082 | 0,37 | 0,082 | 0,434 |
| C15567[0.409] | 0,538 | 0,286 | 0,486 | 0,454 | 0,12 | 0,146 | 0,077 | 0,382 | 0,077 | 0,311 |
| C01263[0.350] | 0,5 | 0,221 | 0,484 | 0,475 | 0,111 | 0,109 | 0,051 | 0,435 | 0,051 | 0,346 |
| C01592[0.133] | 0,469 | 0,366 | 0,343 | 0,3 | 0,126 | 0,144 | 0,136 | 0,23 | 0,136 | 0,221 |
| C08578[0.110] | 0,298 | 0,946 | 0,252 | 0,247 | 0,098 | 0,142 | 0,092 | 0,47 | 0,092 | 0,272 |

$$S_i = \alpha MF_i + (1 - \alpha) \sum_{j=1}^{N} sig(MB_j * Tan_{i,j})$$

Michael Gerlich, Department of Stress and
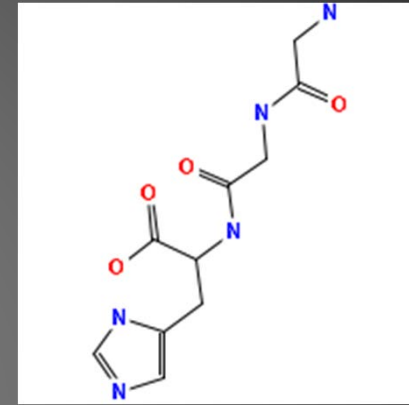Developmental Biology, IPB Halle

# Results

- Test data set: 1099 spectra
  - Secondary metabolites, drugs, toxins, …
  - 344 unique compounds
  - Spectra from Hill et al., RIKEN & IPB
- Median rank of correct compound improved
  - MetFrag: 28
  - MetFusion: 7
- ➢ Works when informative spectra are present, but also when there is loss of information

# MetFrag vs. MetFusion



Cumulative density for MetFrag (black) and MetFusion (red)

Michael Gerlich, Department of Stress and
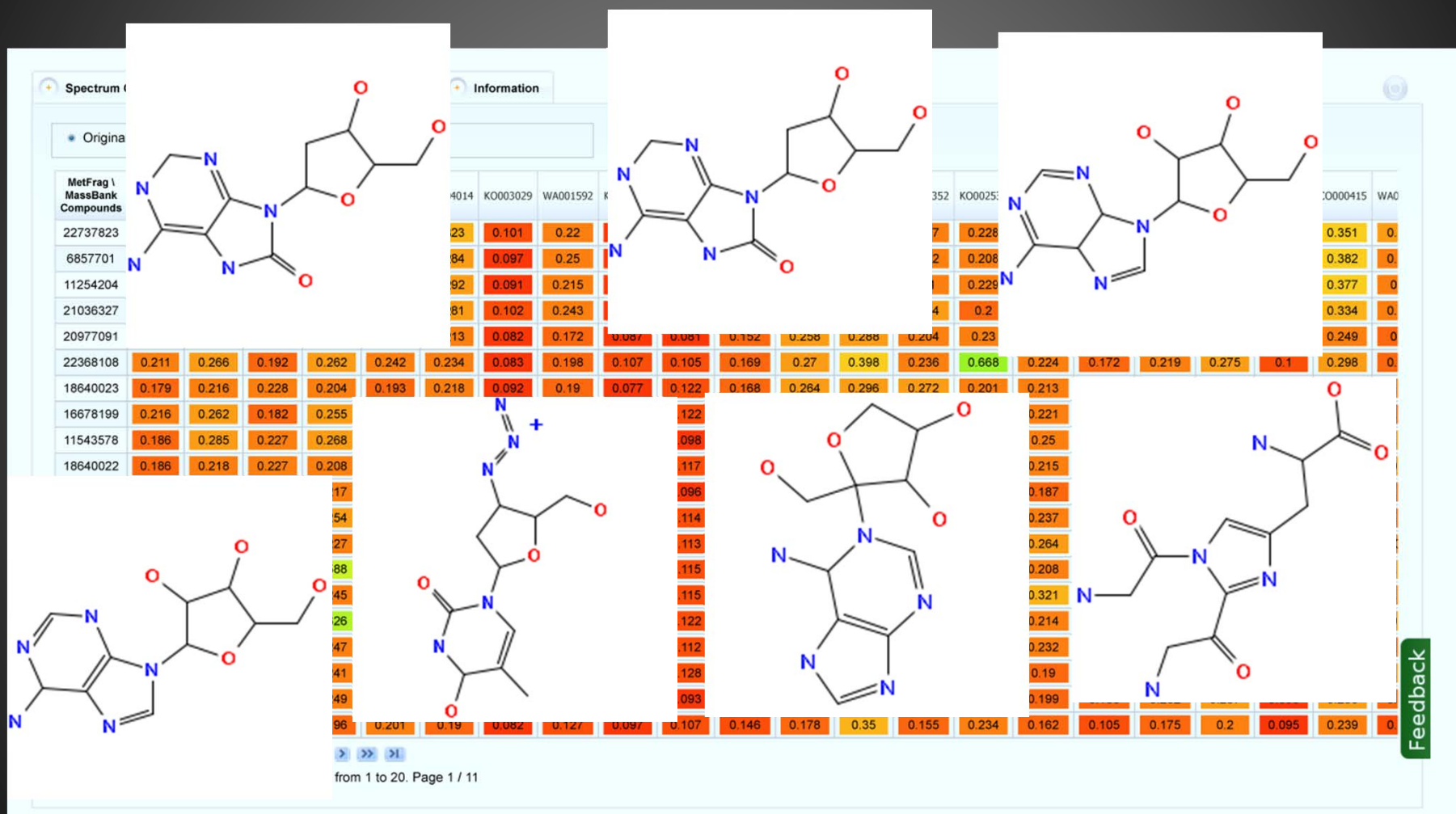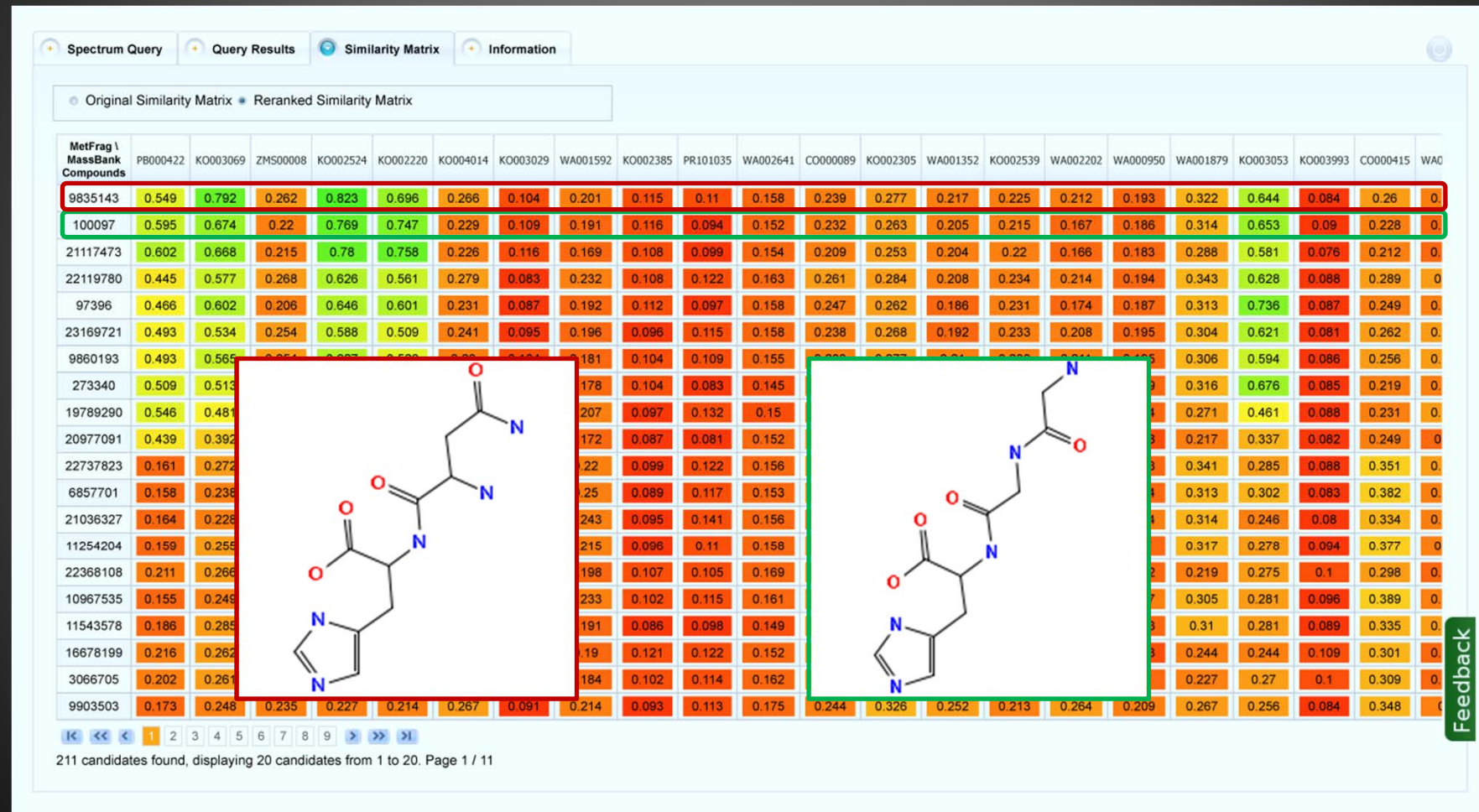Developmental Biology, IPB Halle

# Example

- Tripeptide Gly-Gly-His

- 269.2572 Da

- QqQ spectrum from NIST
  - Nominal masses, modified MetFrag parameters
  - mzabs 0.1 da, mzppm 30ppm


- No tri- or polypeptides present in MassBank
  - But amino acids present

# Original Similarity Matrix

# Re-ranked Similarity Matrix

# Recent additions

- MassBank alternatives
  - Metlin, 54.000 MS2 spectra
  - GMD, 8.800 GC-MS spectra with RI

- HMDB access pending

Michael Gerlich, Department of Stress and
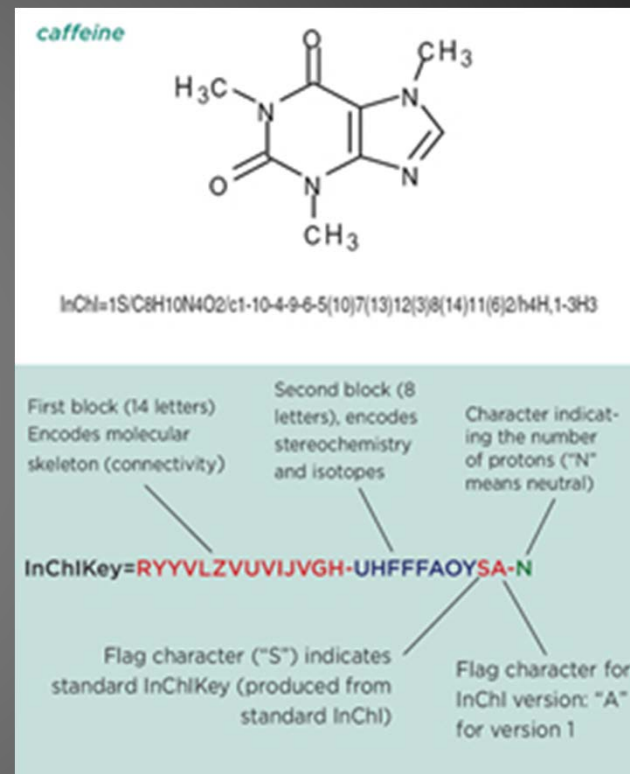Developmental Biology, IPB Halle

# InChIKey-based filtering

- MetFrag candidate list often > 1000
- Lots of stereoisomers per candidate
- Results in clusters, valuable information scarce

- Use connectivity information  from first part of InChIKey to retain only one representative per candidate
- Smaller list of candidates

# InChIKey

- First part connectivity

- Second part stereochemistry



- Image taken from
http://www.iupac.org/publications/ci/2009/3105/iw6_inchi.htm

# Summary

- Combine reference data & prediction

- Improves rank of correct compound

- Access multiple tools within one webpage
  - SDF & XLS export

- Available as web app
  **http://msbi.ipb-halle.de/MetFusion/**

# Acknowledgements

- Dr. Steffen Neumann

- Dr. Emma Schymanski, EAWAG

- Members of MassBank consortium

- Sebastian Wolf, MetFrag developer

- Group members Carsten & Christoph

Thank you for your kind attention