# Metabolite Annotation wit hMetFrag and MetFusion in LC/MS Metabolomics

Steffen Neumann
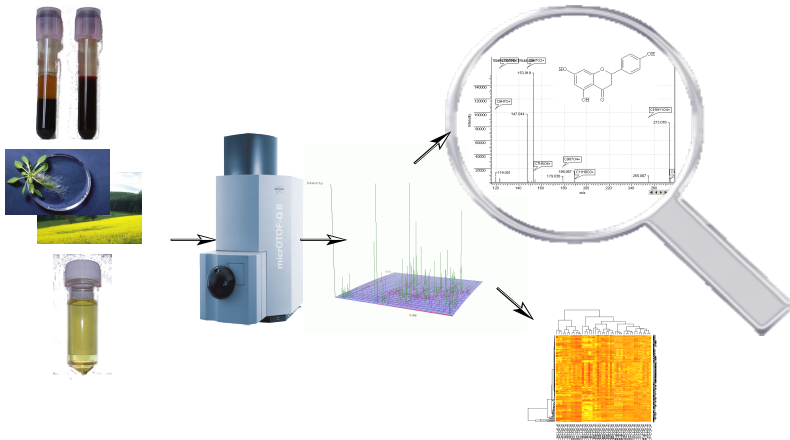
Leibniz Institute of Plant Biochemistry

September, 19$^{th}$ 2014, *NORMAN workshop*

# Leibniz Institute of Plant Biochemistry



- Established 1958 by Kurth Mothes
- Plant diversity, development and adaptation
- Plant production, -protection and biologically active compounds
- About 180 Researchers, including 5-15 Bioinformaticians

Member of the

Leibniz Association

# What *is* Identification ??

"Proposed minimum reporting standards for chemical analysis"
of the Chemical Analysis Working Group (CAWG)
of the Metabolomics Standards Initiative (MSI) for non-novel compounds

1. **Identified compound:**
   retention time/index and mass spectrum, retention time and NMR spectrum, accurate mass and tandem MS, accurate mass and isotope pattern, full $^1$H and/or $^{13}$C NMR, 2-D NMR spectra of in-house measured authentic reference compound Optionally (esp. for unambiguous stereo configuration) selective solvent extraction, retention time, m/z, photodiode array spectra, $\lambda_{max}$ and $\epsilon_{max}$, chemical derivatization, isotope labeling, 2D NMR, IR spectra, etc.

2. **Putative compound:**
   without chemical reference standards, based upon physico-chemical properties and/or spectral similarity with (public/commercial) spectral libraries

3. **Putative compound class**

4. **Unknown compounds** – but can be differentiated and quantified

Member of the
Leibniz Association

# And how confident are you ?

Emma Levels

1. *The* compound structure
2. (Few) Isomers
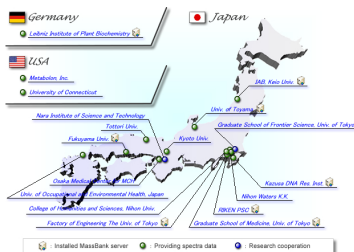3. Compound class
4. Molecular formula
5. Accurate Mass

Schymanski *et al.*, Environmental Science & Technology (2014)

Other (ESI) libraries:

- Metlin
- HMDB
- MMCD
- NIST
- . . .

- http://www.massbank.jp/
- Open Data, Open Consortium
- IPB Halle first European server:
  msbi.ipb-halle.de/MassBank/
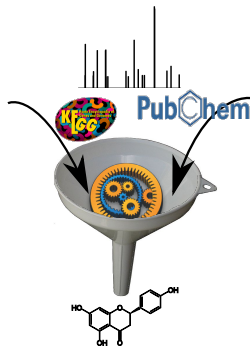- $\approx$ 4 275 compounds
  with $\approx$27 500 MS$^2$ spectra (6/2013)



Horai, H., *et al*. MassBank: a public repository for sharing mass spectral data for life sciences.
Journal of Mass Spectrometry (2010)

# Identification *without* reference spectra

- Spectral libraries (even MassBank) inherently incomplete
- General purpose compound databases:
  - KEGG Compound: 14 067
  - PubChem: 27 million
  - ChemSpider: 25 million
- Known molecular structures
- But: no spectra search $\rightarrow$ "known unknowns"

Wolf, S., Schmidt, S., Müller-Hannemann, M., and Neumann, S.
In silico fragmentation for computer assisted identification of metabolite mass spectra. BMC Bioinf. (2010)

# Identification *without* reference spectra

- Spectral libraries (even MassBank) inherently incomplete
- General purpose compound databases:
  - KEGG Compound: 14 067
  - PubChem: 27 million
  - ChemSpider: 25 million
- Known molecular structures
- But: no spectra search → "known unknowns"

`msbi.ipb-halle.de/MetFrag/` provides this search:

1. Search compound database for precursor mass
2. *In-silico* fragmentation of molecular structure
3. Score measured vs. "predicted" peaks

Wolf, S., Schmidt, S., Müller-Hannemann, M., and Neumann, S.
In silico fragmentation for computer assisted identification of metabolite mass spectra. BMC Bioinf. (2010)

# MetFrag: an example



- Search 290.08Da with 10ppm
- 14 KEGG hits
- MetFrag takes ca. 10sec
- Results ordered by score

- Details / Fragment view
- Excel/SDF download
- Feedback form
- Local version available

Wolf, S., Schmidt, S., Müller-Hannemann, M., and Neumann, S.
In silico fragmentation for computer assisted identification of metabolite mass spectra. BMC Bioinf. (2010)

# MetFrag: an example



- Search 290.08Da with 10ppm
- 14 KEGG hits
- MetFrag takes ca. 10sec
- Results ordered by score

- Details / Fragment view
- Excel/SDF download
- Feedback form
- Local version available

Wolf, S., Schmidt, S., Müller-Hannemann, M., and Neumann, S.
In silico fragmentation for computer assisted identification of metabolite mass spectra. BMC Bioinf. (2010)

## Rest of the world

There are other tools for Metabolite annotation/identification, incomplete list and some literature:
Since before 2010

- ACD Fragmenter (ACD)
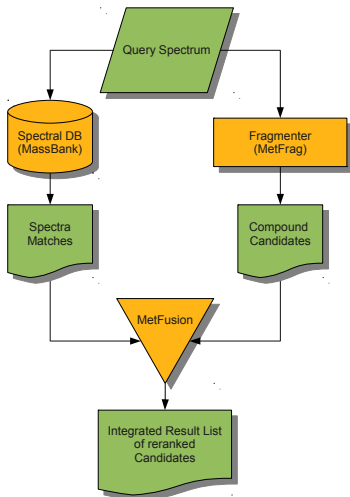- MassFrontier (HighChem), Grant *et al.* Anal. Chem (2008)

After 2010

- MAGMa (Ridder *et al.,*) CASMI 2013
- CFM (Allen *et al.,*) CASMI 2013
- Rational Numbers (Sweeny *et al.,*) CASMI 2013
- SIRIUS (Böcker *et al.,*)
- MassFrontier (HighChem), (Viant *et al.*)

Schymanski, E.L. and Neumann, S. CASMI: And the Winner is . . . Metabolites (2013) / Nishioka *t al.*, Winners of CASMI2013: Automated Tools and Challenge Data, J Mass Spectrometry (2014) / Böcker *et al.*, J. Cheminf

# Integrating MassBank & MetFrag ?

Both tools are great ! Both aid in identifying compounds:

- MassBank: **reference** spectra
    - Actual, real measurements under real conditions
    - Limited coverage, metabolite you look for is (almost) always missing
    - Results based on *spectral* similarity
- MetFrag: **large** compound DB
    - Upstream *exact mass* or *molecular formula* search
    - But: A prediction is a prediction, nothing else
    - Molecular re-arrangements mostly ignored

- MetFusion: **Combine** results:
  - Parallel queries in MassBank and MetFrag
  - Pairwise chemical similarities between result sets
- Calculate *integrated* score
- Major improvement of identification power and accuracy
- `msbi.ipb-halle.de/MetFusion/`

Gerlich M., Neumann S., MetFusion: integration of compound identification strategies
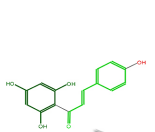Journal of Mass Spectrometry 48 (3), 291-298

**M a s s B a n k**

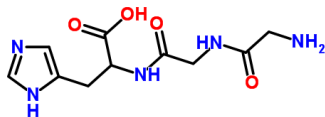| | A | B | C | D | E | F | G | H | I | J | K |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | | C00509[0.975] | C06561[0.965] | C09099[0.956] | C09789[0.916] | C03406[0.599] | C04577[0.520] | C00158[0.502] | C10107[0.468] | C00311[0.418] | -----[0.413] |

# Integrating MassBank + MetFrag: Chemical *Similarity*



**M a s s B a n k**

M
e
t
F
r
a
g

| | A | B | C | D | E | F | G | H | I | J | K |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | | 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] |
| 2 | C00509[1.000] | | | | | | | | | | |
| 3 | C16232[1.000] | | | | | | | | | | |
| 4 | C06561[0.966] | | | | | | | | | | |
| 5 | C12087[0.966] | | | | | | | | | | |
| 6 | C14458[0.966] | | | | | | | | | | |
| 7 | C09826[0.909] | | | | | | | | | | |
| 8 | C03567[0.462] | | | | | | | | | | |
| 9 | C09614[0.462] | | | | | | | | | | |
| 10 | C09751[0.443] | | | | | | | | | | |
| 11 | C09047[0.426] | | | | | | | | | | |
| 12 | C17673[0.426] | | | | | | | | | | |
| 13 | C15567[0.409] | | | | | | | | | | |
| 14 | C01263[0.350] | | | | | | | | | | |
| 15 | C01592[0.133] | | | | | | | | | | |
| 16 | C08578[0.110] | | | | | | | | | | |

Member of the
Leibniz Association

# Integrating MassBank + MetFrag: Chemical *Similarity*



**M a s s B a n k**

| | A | B | C | D | E | F | G | H | I | J | K |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | | 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] |
| 2 | C00509[1.000] | 1,00 | 0,30 | 0,72 | 0,63 | 0,14 | 0,15 | 0,11 | 0,46 | 0,11 | 0,34 |
| 3 | C16232[1.000] | 0,92 | 0,29 | 0,69 | 0,62 | 0,14 | 0,15 | 0,10 | 0,47 | 0,10 | 0,37 |
| 4 | C06561[0.966] | 0,30 | 1,00 | 0,25 | 0,24 | 0,10 | 0,14 | 0,10 | 0,45 | 0,10 | 0,26 |
| 5 | C12087[0.966] | 0,25 | 0,32 | 0,24 | 0,24 | 0,12 | 0,21 | 0,09 | 0,33 | 0,09 | 0,32 |
| 6 | C14458[0.966] | 0,62 | 0,32 | 0,50 | 0,45 | 0,11 | 0,15 | 0,09 | 0,38 | 0,09 | 0,29 |
| 7 | C09826[0.909] | 0,90 | 0,29 | 0,70 | 0,63 | 0,13 | 0,15 | 0,10 | 0,49 | 0,10 | 0,35 |
| 8 | C03567[0.462] | 0,58 | 0,32 | 0,48 | 0,44 | 0,11 | 0,15 | 0,09 | 0,38 | 0,09 | 0,29 |
| 9 | C09614[0.462] | 0,91 | 0,29 | 0,70 | 0,62 | 0,14 | 0,16 | 0,10 | 0,48 | 0,10 | 0,36 |
| 10 | C09751[0.443] | 0,90 | 0,29 | 0,70 | 0,63 | 0,13 | 0,15 | 0,10 | 0,50 | 0,10 | 0,35 |
| 11 | C09047[0.426] | 0,38 | 0,41 | 0,33 | 0,32 | 0,12 | 0,14 | 0,08 | 0,60 | 0,08 | 0,25 |
| 12 | C17673[0.426] | 0,36 | 0,32 | 0,32 | 0,30 | 0,13 | 0,12 | 0,08 | 0,37 | 0,08 | 0,43 |
| 13 | C15567[0.409] | 0,54 | 0,29 | 0,49 | 0,45 | 0,12 | 0,15 | 0,08 | 0,38 | 0,08 | 0,31 |
| 14 | C01263[0.350] | 0,50 | 0,22 | 0,48 | 0,48 | 0,11 | 0,11 | 0,05 | 0,44 | 0,05 | 0,35 |
| 15 | C01592[0.133] | 0,47 | 0,37 | 0,34 | 0,30 | 0,13 | 0,14 | 0,14 | 0,23 | 0,14 | 0,22 |
| 16 | C08578[0.110] | 0,30 | 0,95 | 0,25 | 0,25 | 0,10 | 0,14 | 0,09 | 0,47 | 0,09 | 0,27 |

**M e t F r a g**

Chemical similarity color-coded: high … medium … low

Member of the Leibniz Association

- Example: tripeptide Gly-Gly-His
- Nominal mass spectrum
  from NIST MS/MS library
- PubChem: 211 candidates
- MassBank: only His and Gly spectra
- MetFrag: correct ranked 23$^{rd}$

Gerlich M., Neumann S., MetFusion: integration of compound identification strategies
Journal of Mass Spectrometry 48 (3), 291-298, 2013

# MetFusion: Example Gly-Gly-His



Gerlich M., Neumann S., MetFusion: integration of compound identification strategies
Journal of Mass Spectrometry 48 (3), 291-298, 2013

# MetFusion: Evaluation

- 345 compounds, 89 to 837 Da
- flavonoids, isoflavonoids, steroids, amino acids, carboxylic acids, polyketids, prenol and sterol lipids, glucosides, drugs, toxins, alcohols, carbohydrates
- 1062 spectra from MassBank, all QTOF
- Median 707 candidates from PubChem
- Leave-*some*-out MassBank "pruning":

|  |  | Max. chemical similarity to MassBank |
|---|---|---|
|  |  | =1 |
| MetFusion | Rank | 1 |
|  | RRP | 1 |

Gerlich M., Neumann S., MetFusion: integration of compound identification strategies
Journal of Mass Spectrometry 48 (3), 291-298, 2013

- 345 compounds, 89 to 837 Da
- flavonoids, isoflavonoids, steroids, amino acids, carboxylic acids, polyketids, prenol and sterol lipids, glucosides, drugs, toxins, alcohols, carbohydrates
- 1062 spectra from MassBank, all QTOF
- Median 707 candidates from PubChem
- Leave-*some*-out MassBank "pruning":

|  |  | Max. chemical similarity to MassBank | |
| --- | --- | --- | --- |
|  |  | <1 | =1 |
| MetFusion | Rank | 4 | 1 |
|  | RRP | 0.993 | 1 |

Gerlich M., Neumann S., MetFusion: integration of compound identification strategies
Journal of Mass Spectrometry 48 (3), 291-298, 2013

- 345 compounds, 89 to 837 Da
- flavonoids, isoflavonoids, steroids, amino acids, carboxylic acids, polyketids, prenol and sterol lipids, glucosides, drugs, toxins, alcohols, carbohydrates
- 1062 spectra from MassBank, all QTOF
- Median 707 candidates from PubChem
- Leave-*some*-out MassBank "pruning":

|           |      | Max. chemical similarity to MassBank | | |
| --------- | ---- | ----- | ----- | --- |
|           |      | **<0.9** | <1 | =1 |
| MetFusion | Rank | **7** | 4 | 1 |
|           | RRP  | **0.991** | 0.993 | 1 |

- 345 compounds, 89 to 837 Da
- flavonoids, isoflavonoids, steroids, amino acids, carboxylic acids, polyketids, prenol and sterol lipids, glucosides, drugs, toxins, alcohols, carbohydrates
- 1062 spectra from MassBank, all QTOF
- Median 707 candidates from PubChem
- Leave-*some*-out MassBank "pruning":

|  |  | Max. chemical similarity to MassBank | | | | |
|---|---|---|---|---|---|---|
|  |  | <0.7 | <0.8 | **<0.9** | <1 | =1 |
| MetFusion | Rank | 10 | 8 | **7** | 4 | 1 |
|  | RRP | 0.986 | 0.990 | **0.991** | 0.993 | 1 |

# Outline

# What is CASMI ?

We invited the experimental and computational mass spectrometry community to participate in an open **contest** on the **identification** of small molecules from mass spectrometry data:

## **C**ritical **A**ssessment of **S**mall **M**olecule **I**dentification

www.casmi-contest.org





E. Schymanski, S. Neumann

Isoprothiolane



Challenge5.txt

Isoprothiolane



ES−category2−Challenge5.txt

## Isoprothiolane



Challenge5.txt



**ES−category2−Challenge5.txt**



Score >
- ■ 0.5
- ■ 0.55
- ■ 0.6
- ■ 0.65
- ■ 0.7
- ■ 0.75

Tetrahydroalstonine
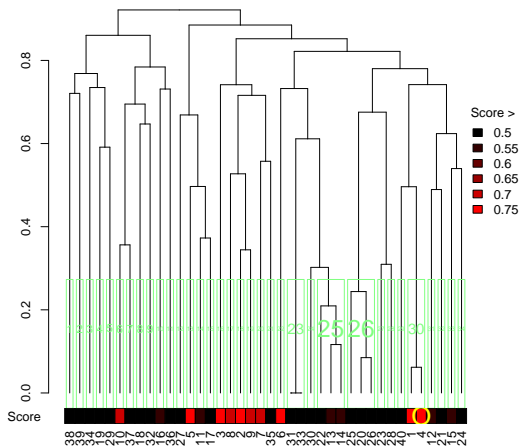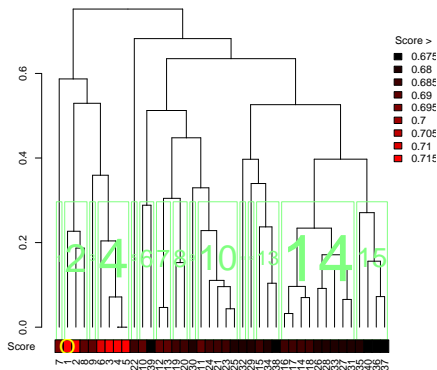
# CASMI as an exercise for post-processing analysis



Acetyl-Gln-Leu-amide

## KEGG extrapolation

The higher the similarity between an unknown and the reference spectrum, the better the identification result.
⇒ How similar are compounds between MassBank and KEGG ?

Number of KEGG compounds for which a MassBank record with a chemical similarity greater or equal the threshold exists:

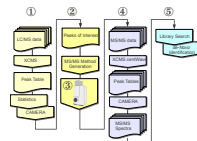| Similarity | >0.7 | >0.8 | **>0.9** | =1.0 |
|---|---|---|---|---|
| KEGG Entries | 5 513 | 4 068 | **2 690** | 1 470 |

MetFusion has median rank 7 on test data if reference spectra with >0.9 chemical similarity are available. **If** that generalises to all KEGG compounds, we expect for 1 345 metabolites rank 7 or better.

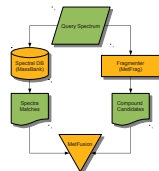Member of the
*Leibniz*
Leibniz Association

## Summary

MetShot aproach

- High-quality MS/MS spectra
- Biologically relevant features
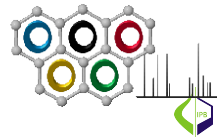- Signal processing with XCMS and CAMERA

Metabolite identification

- With reference spectra (MassBank)
- With *in silico* tools (MetFrag)
- With a combination thereof (MetFusion)

CASMI 2012 – 2013 – 2014

- Open contest for structure elucidation
- Good results, but manual methods don't scale
- Becoming a regular contest series

# Thanks to . . .



MSBI @ IPB Halle:

- Christoph Ruttkies (MetFrag)
- Michael Gerlich (MetFusion)
- D. Schober, S. Mönchgesang

Acknowledgements:

- Profs. Nishioka, Arita (MassBank Japan)
- Emma Schymanski (Eawag, Zürich)
- Funding: DFG grant "ChemFrag",
  EU FP7 "COSMOS",
  EU FP7 "SOLUTIONS"

Alumni (excerpt):

- Dr. Sebastian Wolf (now Bruker Biospin)
- Dr. Ralf Tautenhahn (SCRIPPS, now ThermoFischer)
- Carsten Kuhl (now Bruker Biospin)
- Björn Egert (now Max Rübner Institute)



Member of the
Leibniz Association

Only supplemental slides
beyond this point