*Non-Target Screening:*
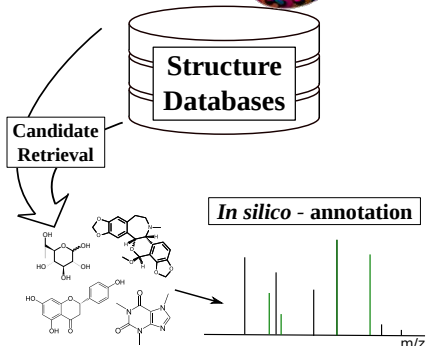*In silico* fragmentation and
reference information

Christoph Ruttkies

Leibniz Institute of Plant Biochemistry
Department of Stress and Developmental Biology
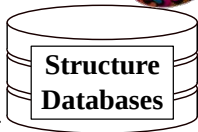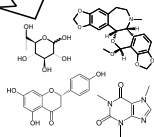Bioinformatics & Mass Spectrometry

16-9-2014

s□luti♥ns
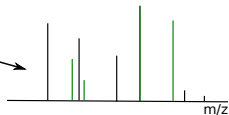
## Workflow

## Workflow

## Workflow



**Candidate
selection**

**Fragmentation**

# In silico fragmentation

## *In silico* fragmentation



- here: combinatorial

- also rule based fragmentation approaches[1]

[1] Kerber, A. *et al.* (2013): Mathematical Chemistry and Chemoinformatics: Structure
Generation, Elucidation and Quantitative Structure-Property Relationships

## Workflow



**Candidate
selection**

**Structure
Databases**

**Candidate
Retrieval**

*In silico* - **annotation**

**Fragmentation**

$S_{C_1} = 1.0$
$S_{C_2} = 0.89$
$S_{C_3} = 0.63$

**Scoring/
Ranking**

## Scoring



- each candidate for a spectrum is scored separately
- → **mass** of matched peak is considered
- → **intensity** of peak is considered
- → **strength of bonds broken** are considered

## Where are we? What do we want to achieve?

**Current state:**

- HR-MS/MS data of molecules from environmental reference standards
- benchmark data of 466 merged spectra (+, -)

- MetFrag ranks 12 % correctly in top spot (PubChem + mol. formula)

**Aim:**

improve HR-MS/MS based identification (MetFrag) with additional information from the experimental context

references in literature (PubMed)

## Where are we? What do we want to achieve?

**Current state:**

- HR-MS/MS data of molecules from environmental reference standards
- benchmark data of 466 merged spectra (+, -)

- MetFrag ranks 12 % correctly in top spot (PubChem + mol. formula)

**Aim:**

improve HR-MS/MS based identification (MetFrag) with additional information from the experimental context

**references in literature (PubMed)**

s**e**luti**e**ns

## Where are we? What do we want to achieve?

**Current state:**

- HR-MS/MS data of molecules from environmental reference standards
- benchmark data of 466 merged spectra (+, -)

- MetFrag ranks 12 % correctly in top spot (PubChem + mol. formula)

**Aim:**

improve HR-MS/MS based identification (MetFrag) with additional information from the experimental context

**references in literature (PubMed) ???**

solutions

## Preliminary work

Little, J.L. *et al.* (2011): **Identification of "Known Unknowns" Utilizing Accurate Mass Data and ChemSpider**[1]

**"Known Unknowns"** - unknown to the investigator but known to the chemical literature, reference database, internet resource

**Approach:**

- selecting compounds by accurate mass search (ChemSpider, PubChem)
- sorting them (descending) by number of references
- most useful results at the top

**The more an environmental compound is used the more likely it is to be found in a sample.**

[1] J. Am. Soc. Mass Spectrom. 23, 179-185

s☐luti☐ns

## Preliminary work

Little, J.L. *et al.* (2011): **Identification of "Known Unknowns" Utilizing Accurate Mass Data and ChemSpider**[1]

**"Known Unknowns"** - unknown to the investigator but known to the chemical literature, reference database, internet resource

**Approach:**

- selecting compounds by accurate mass search (ChemSpider, PubChem)
- sorting them (descending) by number of references
- most useful results at the top

**The more an environmental compound is used
the more likely it is to be found in a sample.**

[1] J. Am. Soc. Mass Spectrom. 23, 179-185

## Preliminary work

Little, J.L. *et al.* (2011): **Identification of "Known Unknowns" Utilizing Accurate Mass Data and ChemSpider**[1]

**"Known Unknowns"** - unknown to the investigator but known to the chemical literature, reference database, internet resource

**Approach:**

- selecting compounds by accurate mass search (ChemSpider, PubChem)
- sorting them (descending) by number of references
- most useful results at the top

no HR-MS/MS info included

**The more an environmental compound is used the more likely it is to be found in a sample.**

## Reference resources for chemicals

# **Scopus**

**Pub Med**

**Google** scholar

**SciFinder®**
Essential content. Proven results.

**ChemSpider**
The free chemical database

## **WEB OF SCIENCE**

**Why we use PubMed?** → **Content**

- technical reasons
- open access
- PubChem (direct link)

- 23M citations (incl. MEDLINE)
- biomedical, chemical, biological, environmental[1]

[1] http://www.nlm.nih.gov/pubs/factsheets/medline.html

**solutions**

## Reference resources for chemicals

# **Scopus**

**Pub Med**

**Google** scholar

SciFinder®
Essential content. Proven results.™

ChemSpider
The free chemical database

### **WEB OF SCIENCE**

**Why we use PubMed?** → Content

- technical reasons
- open access
- PubChem (direct link)

- 23M citations (incl. MEDLINE)
- biomedical, chemical, biological, environmental[1]

[1]http://www.nlm.nih.gov/pubs/factsheets/medline.html

selutions

## Reference resources for chemicals

# Scopus

**Pub Med**

Google scholar

SciFinder®
Essential content. Proven results.™

ChemSpider
The free chemical database

## WEB OF SCIENCE

**Why we use PubMed?** → **Content**

- technical reasons
- open access
- PubChem (direct link)

- 23M citations (incl. MEDLINE)
- biomedical, chemical, biological, environmental[1]

[1] http://www.nlm.nih.gov/pubs/factsheets/medline.html

solutions

# Example: Chlorpyriphos's PubMed References

# Example: Chlorpyriphos's PubMed References

## Combining MetFrag with References

**Pub🔷hem**

Structure
Database

## Combining MetFrag with References



Candidates

## Combining MetFrag with References

## Combining MetFrag with References

## Combining MetFrag with References



$\alpha \cdot \mathrm{MetFrag}_C$      $+$      $\beta \cdot \mathrm{nPubMedRefs}_C$

## Combining MetFrag with References



$\alpha \cdot \mathrm{MetFrag}_C$ $+$ $\beta \cdot \mathrm{nPubMedRefs}_C$

$$S_{C_1} = 1.0$$

$$S_{C_2} = 0.97$$

$$S_{C_3} = 0.71$$
$$\vdots$$

## Example: Chlorpyriphos

HR-MS/MS data of Challenge 9 (Chlorpyrifos)
CASMI contest 2013[1]

[1] http://casmi-contest.org/2013/challenges.shtml

## Example: Chlorpyriphos

HR-MS/MS data of Challenge 9 (Chlorpyrifos)
CASMI contest 2013[1]

$\rightarrow$ PubChem + 348.9264 Da = 114 candidates

$\rightarrow$ MetFrag alone: Rank 3

# Example: Chlorpyriphos

HR-MS/MS data of Challenge 9 (Chlorpyrifos)
CASMI contest 2013[1]

$\rightarrow$ PubChem + 348.9264 Da = 114 candidates

$\rightarrow$ MetFrag alone: Rank 3



(1) CID: 57354037  (2) CID: 13274485  (3) Chlorpyrifos

---

[1]http://casmi-contest.org/2013/challenges.shtml

## Example: Chlorpyriphos

HR-MS/MS data of Challenge 9 (Chlorpyrifos)
CASMI contest 2013[1]

$\rightarrow$ PubChem + 348.9264 Da = 114 candidates

$\rightarrow$ MetFrag alone: Rank 3



(1) CID: 57354037   (2) CID: 13274485   (3) Chlorpyrifos
No. Refs: 0             No. Refs: 0             No. Refs: 151

[1] http://casmi-contest.org/2013/challenges.shtml

## Example: Chlorpyriphos

HR-MS/MS data of Challenge 9 (Chlorpyrifos)
CASMI contest 2013[1]

→ PubChem + 348.9264 Da = 114 candidates

→ MetFrag + PubMed Refs: Rank 1



| (1) CID: 57354037 | (2) CID: 13274485 | (3) Chlorpyrifos |
| No. Refs: 0 | No. Refs: 0 | No. Refs: 151 |

[1] http://casmi-contest.org/2013/challenges.shtml

## Number of ranked first correctly (466 spectra)



$$S_c = \alpha \cdot \mathrm{MetFrag}_c + \beta \cdot \mathrm{nPubMedRefs}_c$$

## Number of ranked first correctly (466 spectra)



$$S_C = \alpha \cdot \text{MetFrag}_C + \beta \cdot \text{nPubMedRefs}_C$$

## Number of ranked first correctly (466 spectra)



$$S_C = \alpha \cdot \mathrm{MetFrag}_C + \beta \cdot \mathrm{nPubMedRefs}_C$$

## Aim: Further improvement

- include further parameters to improve identification
- information weighted based on (user-)defined priorities

$$S_C = \omega_1 \cdot (\textbf{MetFrag})_C$$
$$+ \omega_2 \cdot (\textbf{PubMedRefs})_C$$

s☐luti☐ns

## Aim: Further improvement

- include further parameters to improve identification
- information weighted based on (user-)defined priorities

$$
\begin{aligned}
S_C = \ &\omega_1 \cdot (\textbf{MetFrag})_C \\
+ \ &\omega_2 \cdot (\textbf{PubMedRefs})_C
\end{aligned}
$$

first ranked

$\sim 12\ \%$

$\sim 65\ \%$

s⚫luti⚫ns

## Aim: Further improvement

- include further parameters to improve identification
- information weighted based on (user-)defined priorities

$$
\begin{aligned}
S_C = \ &\omega_1 \cdot (\textbf{MetFrag})_C & \sim 12\ \% \\
+\ &\omega_2 \cdot (\textbf{PubMedRefs})_C & \sim 65\ \% \\
+\ &\omega_3 \cdot \text{Retention time prediction} &
\end{aligned}
$$

first ranked

s●luti♥ns

## Retention time modelling



$$S_C = \alpha \cdot \mathrm{MetFrag}_C + \beta \cdot \mathrm{nPubMedRefs}_C + \gamma \cdot \mathrm{RT}_C$$

s●luti●ns

## Aim: Further improvement

- include further parameters to improve identification
- information weighted based on (user-)defined priorities

$$
\begin{aligned}
S_C = \ & \omega_1 \cdot (\textbf{MetFrag})_C \\
& + \omega_2 \cdot (\textbf{PubMedRefs})_C \\
& + \omega_3 \cdot \text{Retention time prediction}
\end{aligned}
$$

first ranked

$\sim 12\ \%$

$\sim 65\ \%$

$\sim 66\ \%$

s**₀**luti**♥**ns

## Cases with no references

For **134** of 466 spectra the correct candidate has no PubMed reference.



$$S_C = \alpha \cdot \text{MetFrag}_C + \beta \cdot \text{nPubMedRefs}_C$$

## Cases with no references

For **134** of 466 spectra the correct candidate has no
PubMed reference.



With retention time information
we got 35 (26 %) ranked first.

$$S_C = \alpha \cdot \mathrm{MetFrag}_C + \beta \cdot \mathrm{RT}_C$$

## Aim: Further improvement

- include further parameters to improve identification
- information weighted based on (user-)defined priorities

$$
\begin{aligned}
S_C = \ & \omega_1 \cdot (\textbf{MetFrag})_C && \sim 12\ \% \\
+ \ & \omega_2 \cdot (\textbf{PubMedRefs})_C && \sim 65\ \% \\
+ \ & \omega_3 \cdot \text{Retention time prediction} && \sim 66\ \% \\
+ \ & \omega_4 \cdot \text{Spectral information} && ? \\
+ \ & \omega_5 \cdot \text{Substructure (non$-$)presence} && ? \\
+ \ & \omega_6 \cdot \text{Effect prediction} && ? \\
+ \ & \omega_i \cdot ...
\end{aligned}
$$

first ranked

**Future aim:** create an interface that allows user-defined scoring terms

s●luti●ns

## Acknowledgement

Dierk Scheel            IPB Halle
Steffen Neumann         IPB Halle
Emma Schymanski         Eawag Zurich
Carsten Kuhl            IPB Halle
Michael Gerlich         IPB Halle
Susann Mönchgesang      IPB Halle
Heinz Singer            Eawag Zurich
Michael Stravs          Eawag Zurich
Juliane Hollender       Eawag Zurich

## Acknowledgement

Dierk Scheel          IPB Halle
Steffen Neumann       IPB Halle
Emma Schymanski       Eawag Zurich
Carsten Kuhl          IPB Halle
Michael Gerlich       IPB Halle
Susann Mönchgesang    IPB Halle
Heinz Singer          Eawag Zurich
Michael Stravs        Eawag Zurich
Juliane Hollender     Eawag Zurich

**Thank you for attention!**

Q&A?

# Appendix

## MetFrag with References: Results

- 333 of 421 MS/MS spectra have candidates with PubMed entries

- **Question:** Do we need MetFrag when there is a PubMed entry?

- **Answer:** Show results on the 333 spectra where PubMed entries are present in the candidate set.

solutions

## MetFrag with References: Results

MS/MS dataset of 333 merged spectra



$$S_C = \alpha \cdot \mathrm{MetFrag}_C + (1 - \alpha) \cdot \mathrm{nPubMedRefs}_C$$

## MetFrag with References: Results

MS/MS dataset of 333 merged spectra



$$S_C = \alpha \cdot \mathrm{MetFrag}_C + (1 - \alpha) \cdot \mathrm{nPubMedRefs}_C$$

## CASMI 2013 - Challenge 9



(a) 57354037          (b) 13274485          (c) Chlorpyrifos

## User-defined scoring terms

- user provides a list of structures with determined features

| CID | Toxicity | Effect $x_1$ | Effect $x_2$ | Pred.Ret.time | ... |
|---|---|---|---|---|---|
| 321 | 1 | 1 | -1 | 5.63 | ... |
| 26881 | 1 | -1 | -1 | 10.22 | ... |
| 9752 | -1 | -1 | 1 | 11.33 | ... |
| 909 | 1 | 1 | -1 | 9.87 | ... |
| 87665 | 1 | -1 | -1 | 9.41 | ... |

- identification of compounds relies on the given information with user-defined priorities

solutions

# CASMI 2013: Challenge 9

**CASMI** 2013

Critical Assessment of Small Molecule Identification

- MS data for 16 challenges provided
- measured compounds are unknown to the investigator
- use your desired workflow and identification method

- we tackled **Challenge 9** originating from an environmental sample together with **PubMed**

## Normalised PubMed reference score

$$\mathrm{nPubMedRefs}_C = \frac{\mathrm{PubMedRefs}_C}{max(\mathrm{PubMedRefs}_C)}$$

## MetFrag score

score of one candidate c:

$$S_c = \sum_{f \in F_c} \frac{\left(10.0 \cdot \frac{\text{Mass}_f}{\text{Mass}_c}\right)^\alpha \cdot \left(\text{RelInt}_f\right)^\beta}{\left(\sum_{b \in f} \text{BDE}_b\right)^\gamma}$$

## MetFrag score

score of one candidate c:

$$S_c = \sum_{f \in F_c} \frac{\left(10.0 \cdot \frac{\text{Mass}_f}{\text{Mass}_c}\right)^{\alpha} \cdot \left(\text{RelInt}_f\right)^{\beta}}{\left(\sum_{b \in f} \text{BDE}_b\right)^{\gamma}}$$

- all matched fragments f of candidate c

MetFrag score

score of one candidate c:

$$S_c = \sum_{f \in F_c} \frac{\left(10.0 \cdot \frac{\text{Mass}_f}{\text{Mass}_c}\right)^\alpha \cdot (\text{RelInt}_f)^\beta}{\left(\sum_{b \in f} \text{BDE}_b\right)^\gamma}$$

- all matched fragments f of candidate c
- relative mass of matched fragment f and the candidate c

## MetFrag score

score of one candidate c:

$$S_c = \sum_{f \in F_c} \frac{\left(10.0 \cdot \frac{\mathrm{Mass_f}}{\mathrm{Mass_c}}\right)^\alpha \cdot \left(\mathrm{RelInt_f}\right)^\beta}{\left(\sum_{b \in f} \mathrm{BDE_b}\right)^\gamma}$$

- all matched fragments f of candidate c
- relative mass of matched fragment f and the candidate c
- relative intensity of explained peak

## MetFrag score

score of one candidate c:

$$S_c = \sum_{f \in F_c} \frac{\left(10.0 \cdot \frac{\text{Mass}_f}{\text{Mass}_c}\right)^\alpha \cdot \left(\text{RelInt}_f\right)^\beta}{\left(\sum_{b \in f} \text{BDE}_b\right)^\gamma}$$

- all matched fragments f of candidate c
- relative mass of matched fragment f and the candidate c
- relative intensity of explained peak
- sum of Bond Dissociation Energies (BDEs) of all broken bonds b