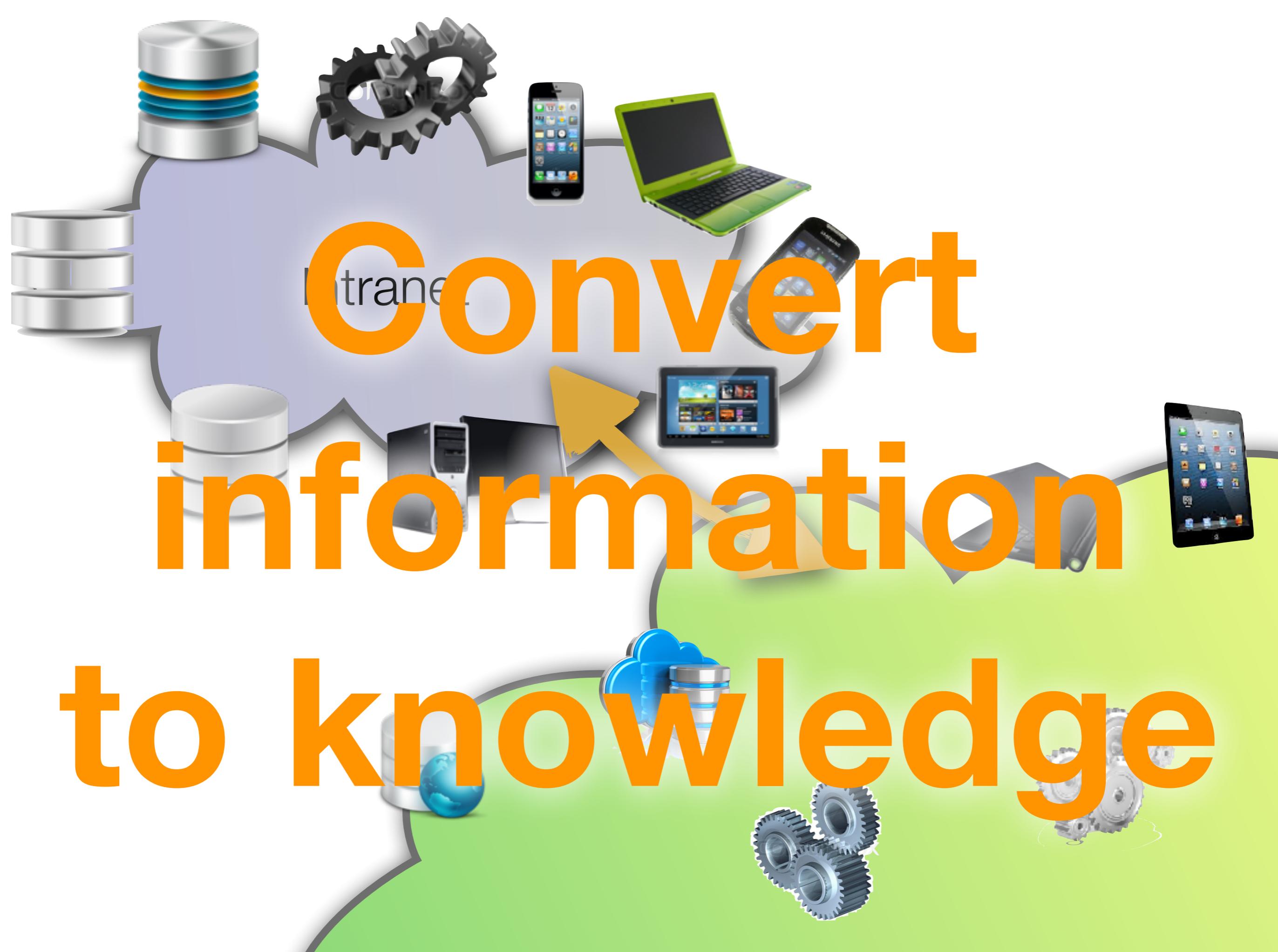




# New challenges in the management of chemical information

---

Service of cheminformatics  
luc.patiny@epfl.ch



tranee  
**Convert**

**information**

**to knowledge**

**bp**

**LD<sub>50</sub>**

**pKa**

**mp**

**IC<sub>50</sub>**

**information**

**NMR**

**IR**

**logP**

**ms**

**Xray**

**QSAR**

**physical property**

**bond length**

**spectra**

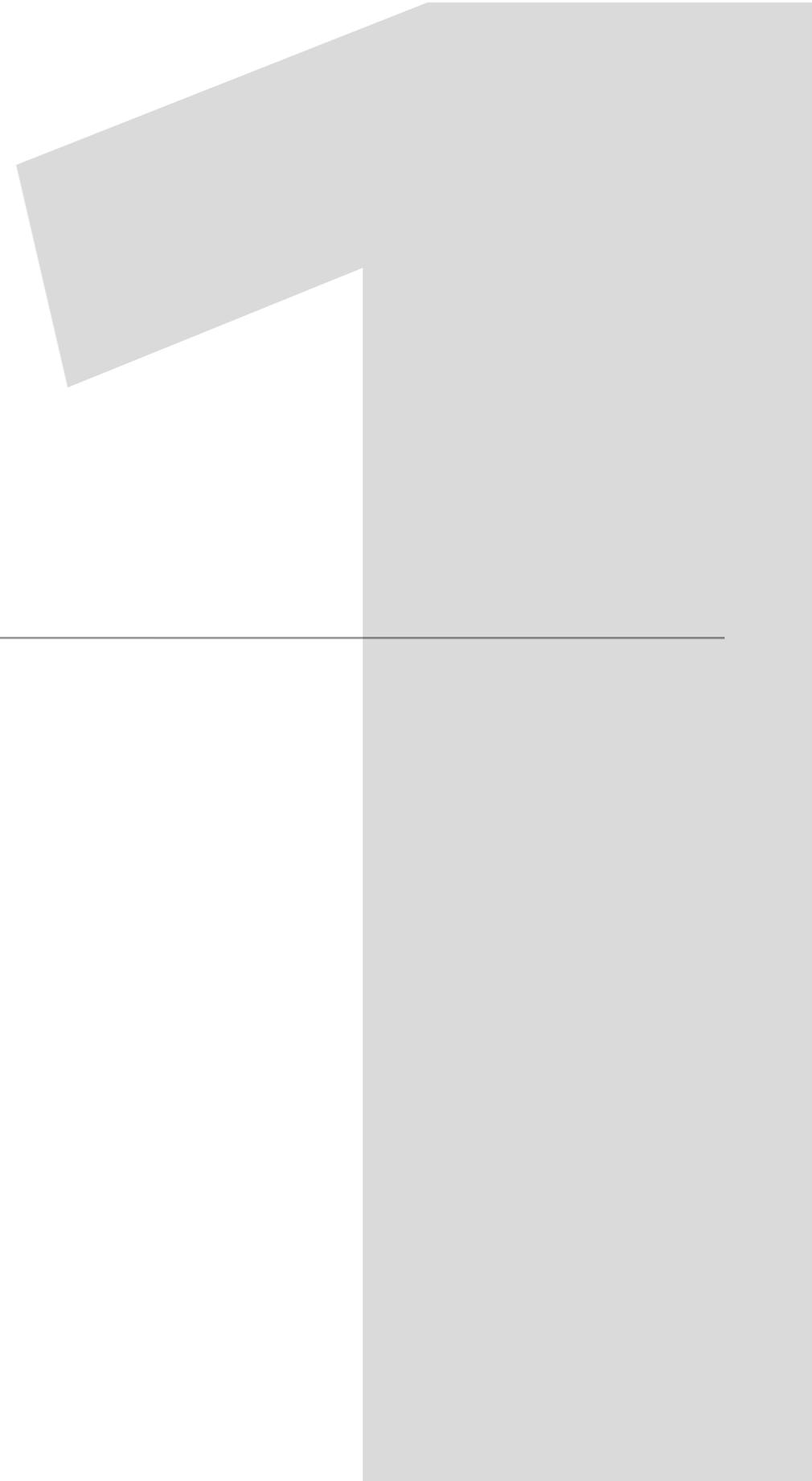
**toxicity**

**fragmentation**

**knowledge**

**Data**

---



# Some available databases

---

ChemSpider

ChemExper

eMolecules

Zinc

# Specialized databases

---

CSD (Cambridge structural database)

KNApSAcK

AIST: Spectral Database

NIST Chemistry web book

# Downloadable databases

---

Crystallography Open Database

DrugBank

Pubchem

GDB-17

ChemBL

PDB

NMR shift DB

KEGG

Mass bank

Binding DB

Wikipedia

OChem : Online chemical database



in 5 years ?

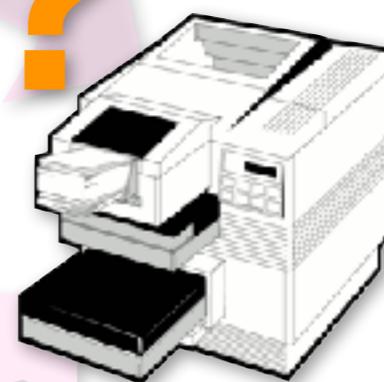


in 10 years ?

FTP  
USB



File Server

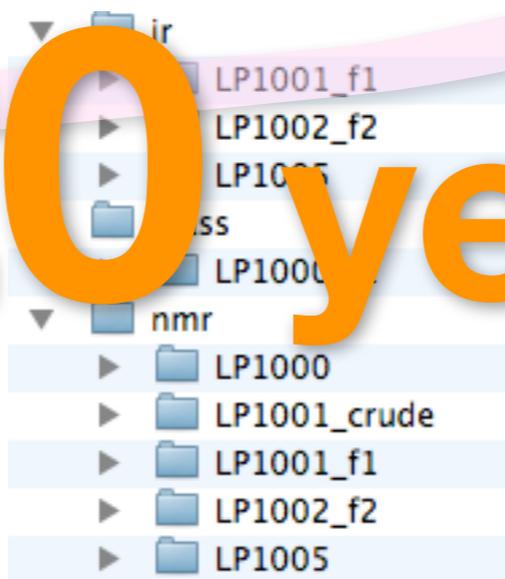


Print  
PDF



in 20 years ?

Install proprietary software for each data format



# Store “correctly” the chemical information

---

- Perennial and reprocessable !
- Only use **ALIVE** published **DESCRIBED** format
  - no proprietary format
  - no PDF, word, excel, ...
  - chemical structures: molfile
  - spectra: JCAMP-DX

# Store “correctly” the chemical information

---

- bp : sign - low - high - pressure
  - 100-120 (20 torr) ; 10 (15 mmHg) ; >300 ; ...
- **NMR** :  $^1\text{H}$  NMR ( $\text{CDCl}_3$ ):  $\delta = 0.84$  (3 H, t,  $J = 7.4$  Hz,  $\text{CH}_3$ ),  $0.94$  (3 H, t,  $J = 7.4$  Hz,  $\text{CH}_3$ ),  $1.23$  [6 H, q,  $J = 6.9$  Hz,  $\text{P}(\text{O})\text{OCH}_2\text{CH}_3$ ],  $1.51$  (4 H, m,  $\text{CH}_2$ ),  $2.20$  (1 H, sextet,  $J = 6.6$  Hz, CH),  $3.80$  (3 H, s,  $\text{OCH}_3$ ),  $4.01$  [4 H, m,  $\text{P}(\text{O})\text{OCH}_2\text{CH}_3$ ],  $4.63$  (1 H, d,  $J_{\text{H,P}} = 17.1$  Hz, NCHP),  $6.88$  (2 H, d,  $J = 8.6$  Hz, CH).
- Traceability, flexibility : noSQL database

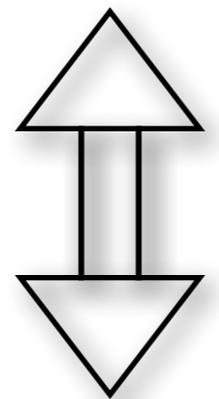
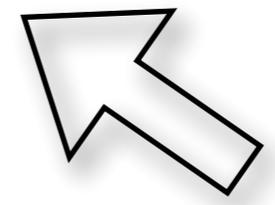
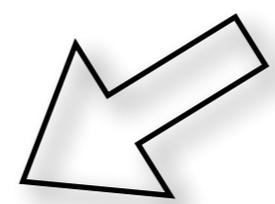
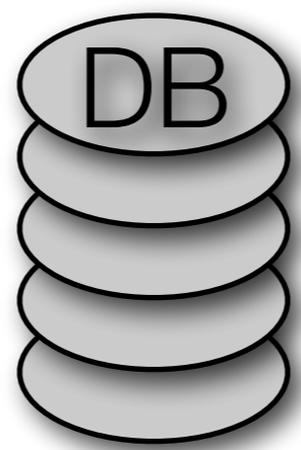
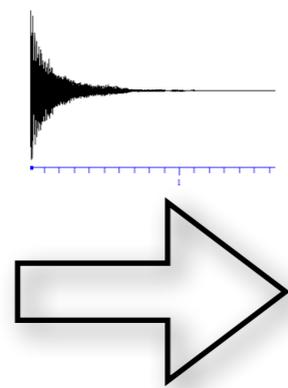
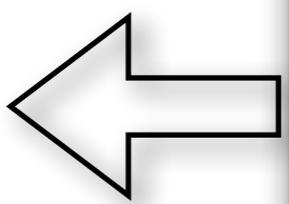


Image analysis  
Mass analysis  
Similarity search  
Predictions

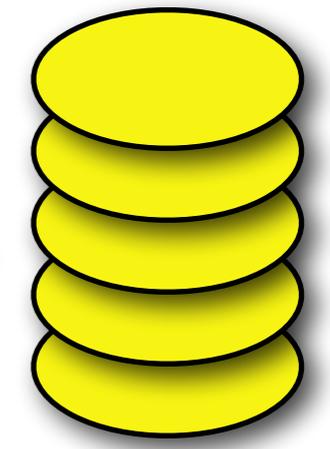
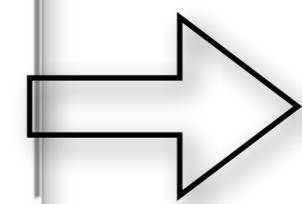


Search for samples

429

Sample ID	Date	Chemical Structure	Mass
178	2013-11-18	<chem>C1=CC=C(C=C1)O</chem>	180.0426
179	2013-11-18	<chem>C1=CC=C(C=C1)O</chem>	180.0426
180	2013-11-18	<chem>C1=CC=C(C=C1)O</chem>	180.0426
181	2013-11-18	<chem>C1=CC=C(C=C1)O</chem>	180.0426

open data

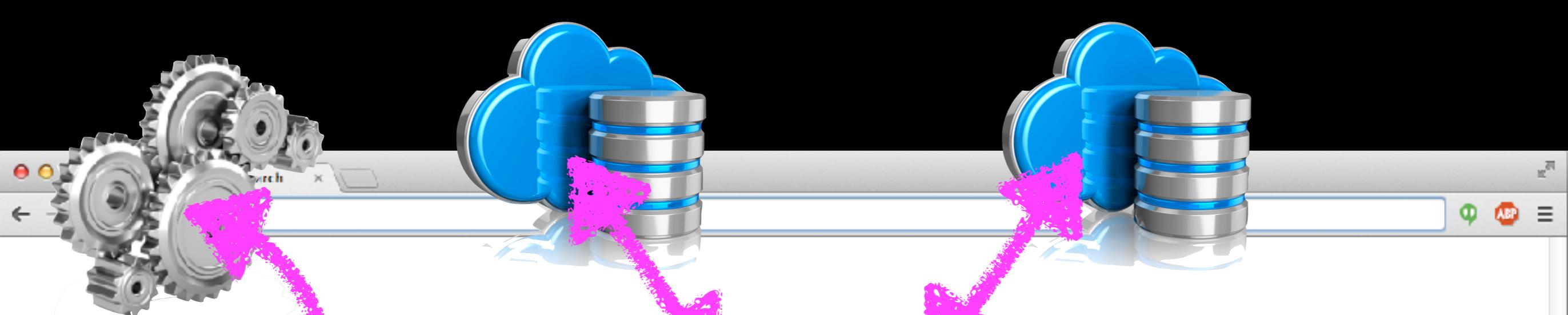


open data

# Retrieving and processing data

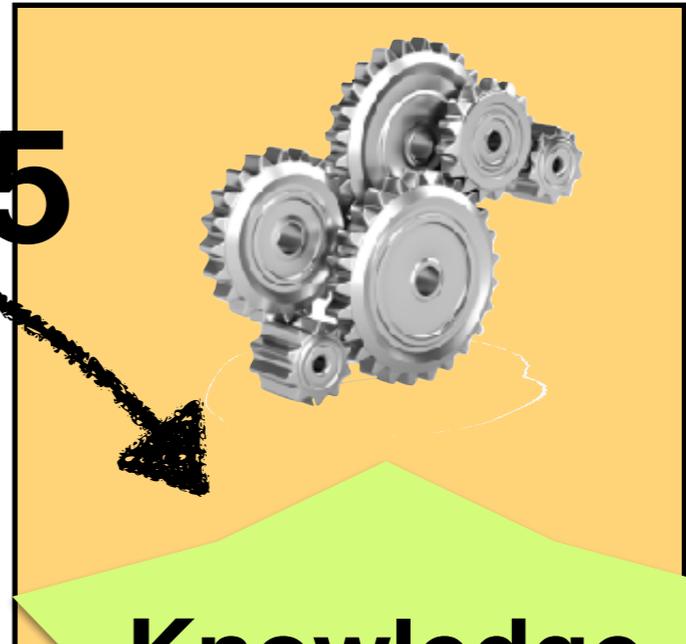
---

The need for a flexible tool



• MacOS, windows, linux

# Pure HTML5



• No installation

• open-source (GitHub)

• No licence fee

Knowledge

Teaching

Research

Unid module

3D scatter plot

Protein sequence

Format: CODATA

ENTRY P918283

SEQUENCE

1 M E T L C Q R L N V C Q D K I L T H Y E N D S T D L R D H I D Y W K H  
E M G F K H I N H Q V V P T L A V S K N K A L  
Y N S O Y S N E K W T L O D V S L E V Y L T A

Un

List

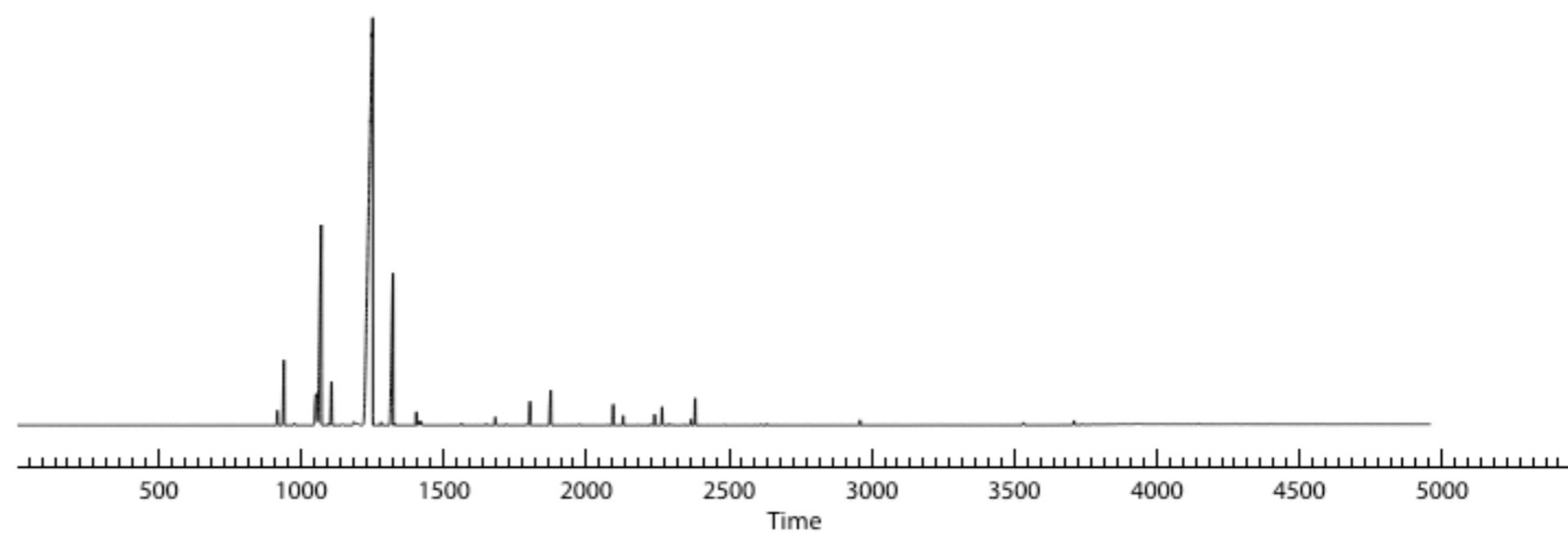
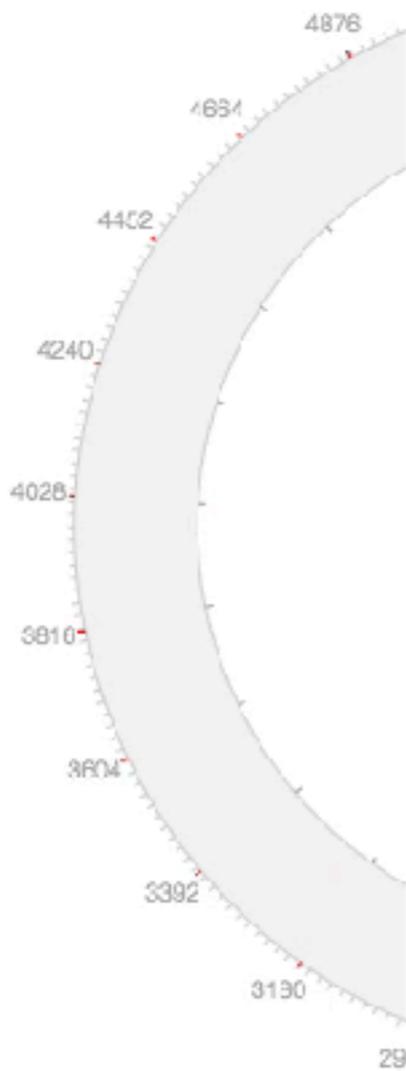
ID

6353

2900

1157

2475



Intensity (-) x 10<sup>-324</sup>

m/z

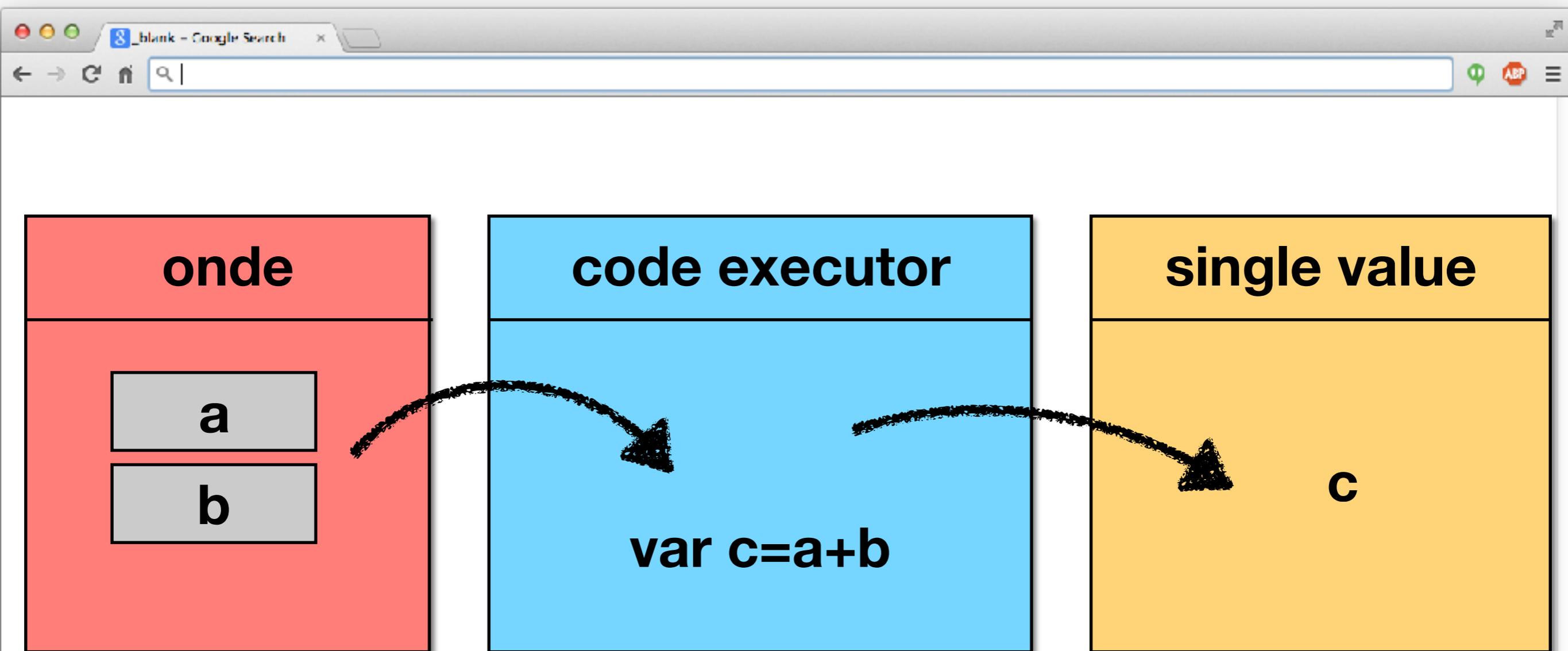
**How does it works ?**

---



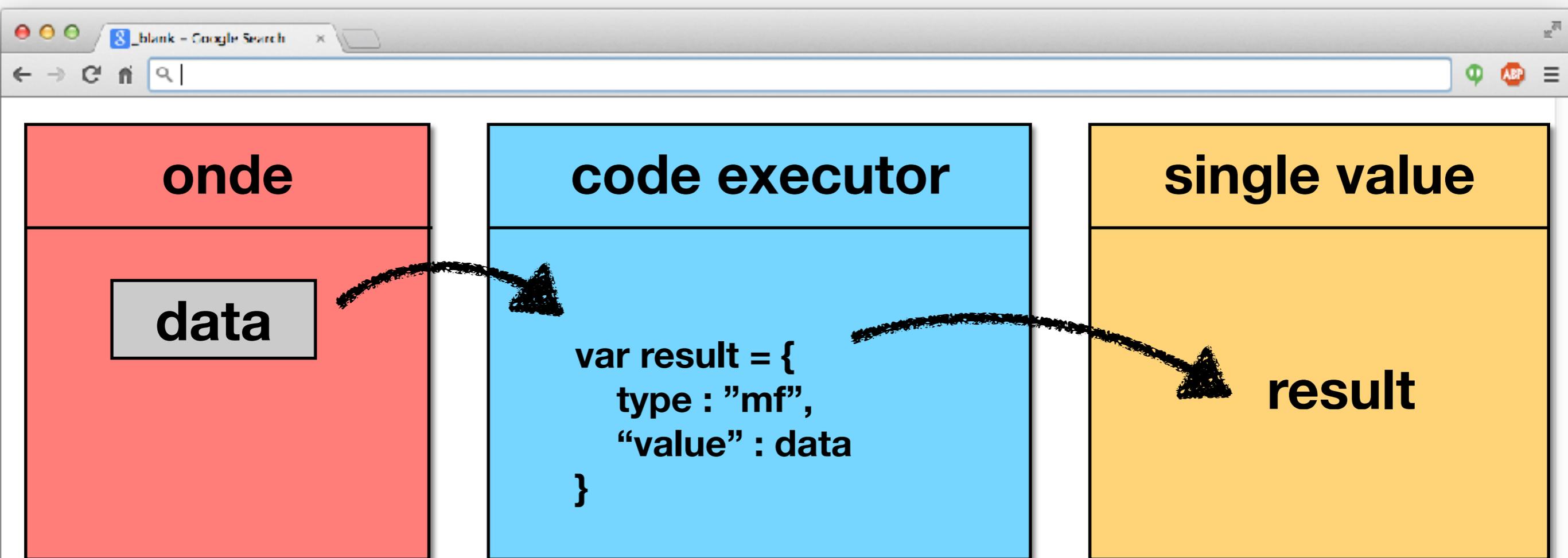
Basic example : calculate a sum of 2 values

---



Demo

# Chemical types - MF and SMILES

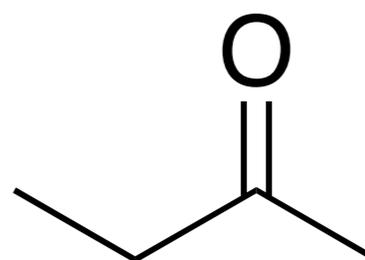


**C6H15N**



**C<sub>6</sub>H<sub>15</sub>N**

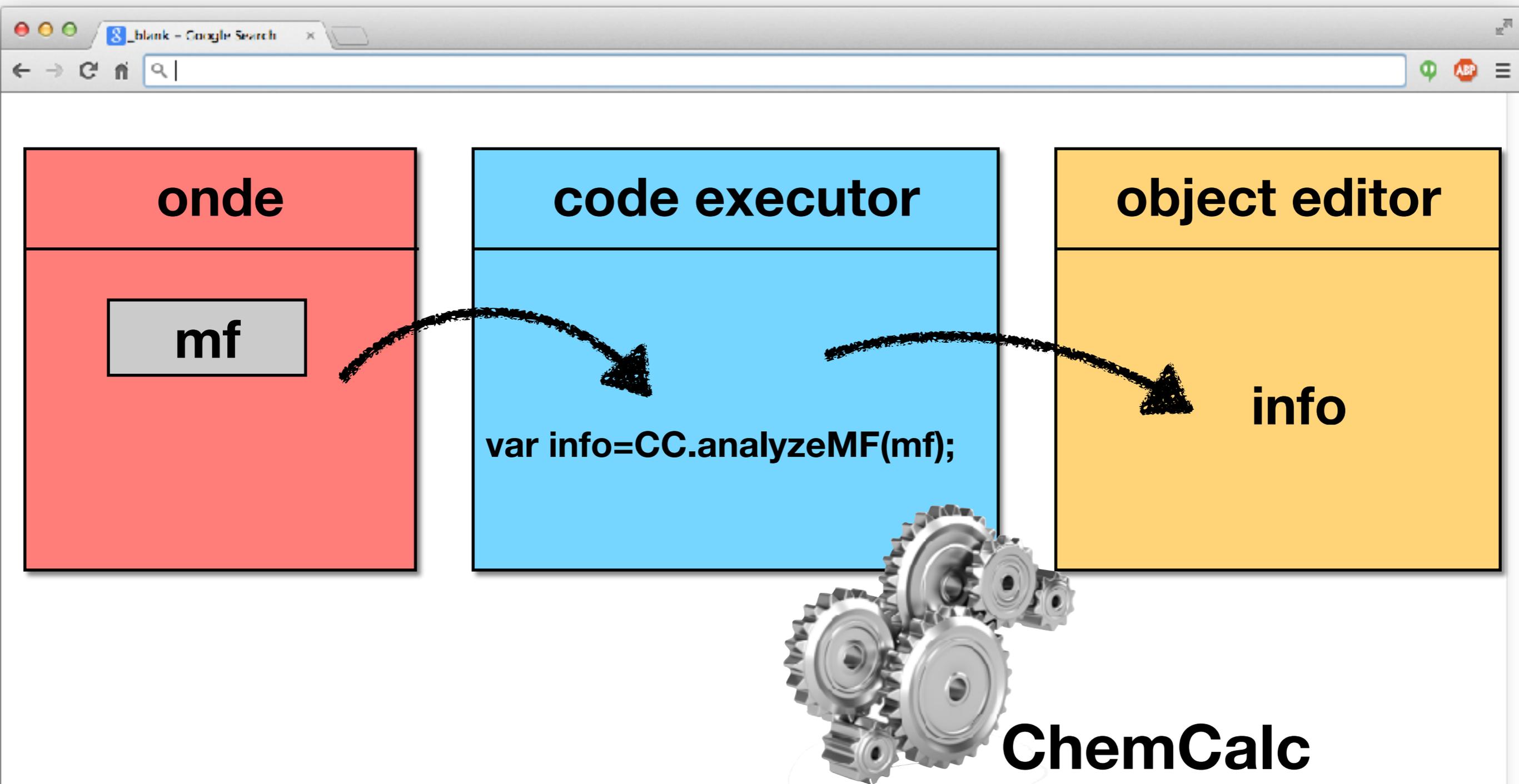
**CCC(=O)C**



Demo

# Using external libraries

---



Demo

# Molecular formula and mass

---

<http://ms.cheminfo.org>

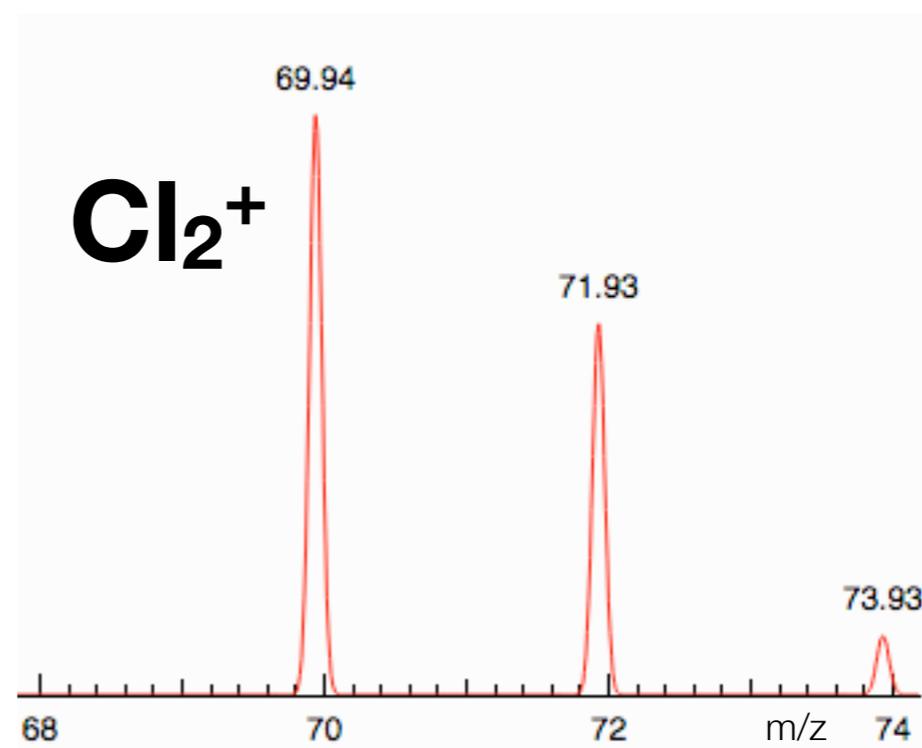
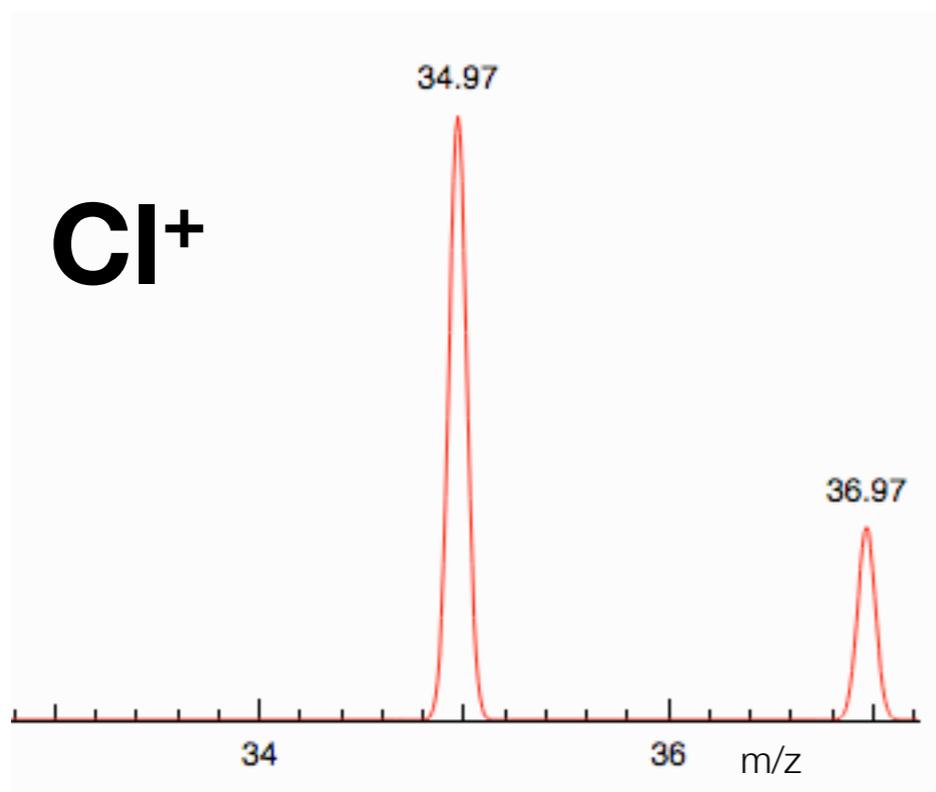
<http://www.chemcalc.org>

# Isotopic distribution

17  
**Cl**  
Chlorine  
35.45

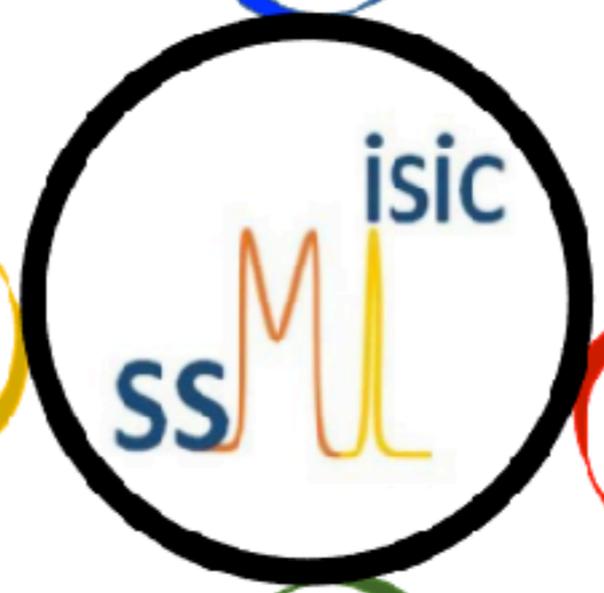
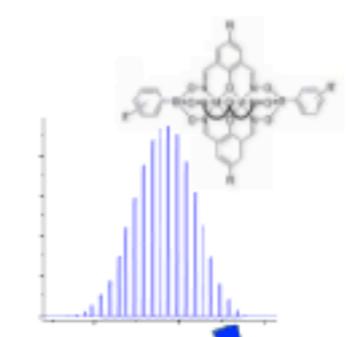
2  
8  
7

Isotope	mass (Da)	Abundance
$^{35}\text{Cl}$	34.97	75.8%
$^{37}\text{Cl}$	36.97	24.2%





Please cite US



Calculate your theoretical spectrum from a molecular formula  
 ↳  $C_6H_6$

From a chemical structure  
 ↳ c1ccc2cc(C(=O)O)ccc2c1

Find a MF from a Monoisotopic Mass  
 ↳ 78.047 →  $C_6H_6$

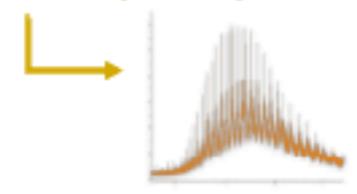
MS/MS fragment calculation  
 ↳ CC(C)C(=O)N → CC(C)C(=O)N + CC(C)C(=O)N

Predict your CID/ETD spectra

↳ Ala-Gly-Ala-Trp

Peptides

Match your Experimental data



Match to list

Drag&Drop and match to contaminants

↳ Easycont

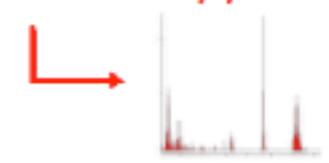
Try our teaching tools

Similarity to spectrum

Extend your search using all elements

S	C	N
O	H	P
Na	K	Cl
Br	I	Ca

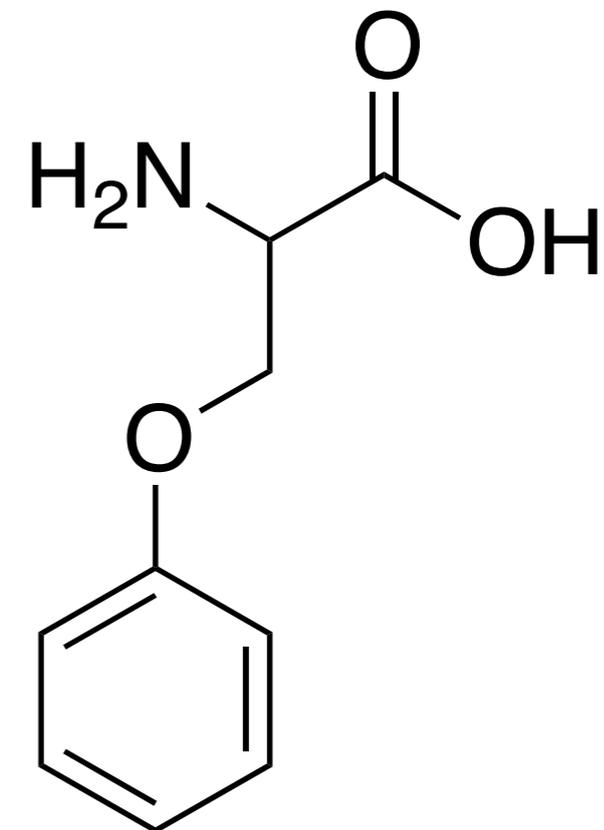
Drag&Drop your Experimental spectrum and Identify your molecular Formula



# ChemCalc

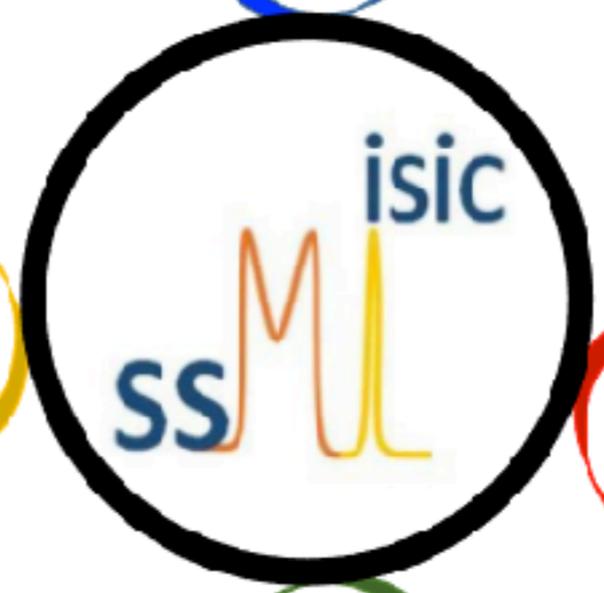
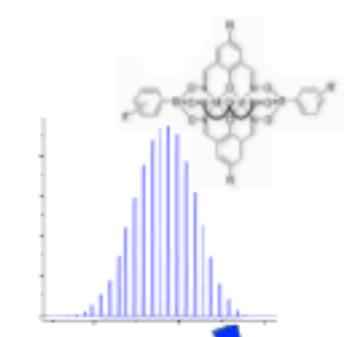
---

- Use of groups (Ph, Ala, Gly, Fmoc)
  - HAlaGlyProOH
- Many components (separated by “.”)
  - C<sub>100</sub> . C<sub>110</sub> . C<sub>120</sub>
- Specific isotope
  - [<sup>13</sup>C]<sub>100</sub> . [<sup>13</sup>C]<sub>50</sub>[<sup>12</sup>C]<sub>50</sub>
- Modification of the isotopic abundance
  - C{50,50}<sub>10</sub>C<sub>10</sub>
- Specification of the charge
  - HAla<sub>10</sub>OH<sup>+</sup> . HAla<sub>10</sub>OH<sup>++</sup> . HAla<sub>10</sub>OH (H<sup>+</sup>)<sub>2</sub>
- Modification of a molecular formula
  - HSer(H<sub>-1</sub>Ph)OH

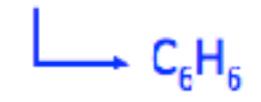




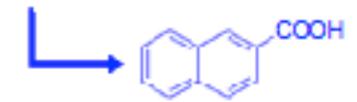
Please cite US



Calculate your theoretical spectrum from a molecular formula



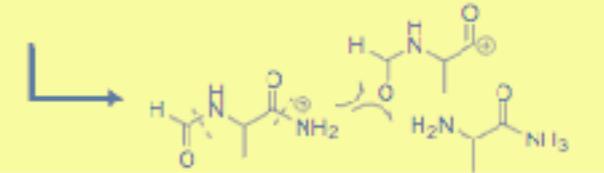
From a chemical structure



Find a MF from a Monoisotopic Mass



MS/MS fragment calculation

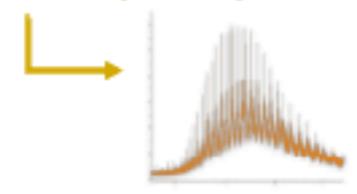


Predict your CID/ETD spectra



Peptides

Match your Experimental data



Try our teaching tools

Match to list

Drag&Drop and match to contaminants

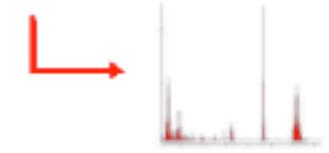


Similarity to spectrum

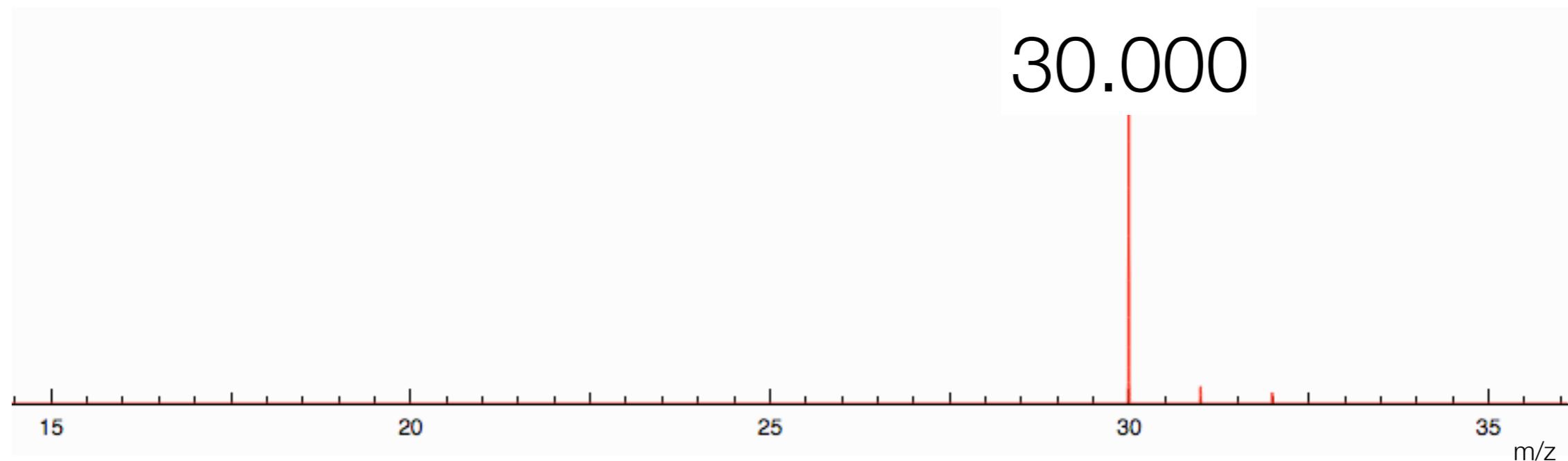
Extend your search using all elements

S	C	N
O	H	F
Br	Cl	I
Na	K	Ca
Fe	Mn	Zn
Al	P	As
Se	Ag	Cd
Cu	Ba	Pb
Li	Be	B
Mg	Si	Ge
Ti	V	Ni
Co	Mo	Sn
Cs	Rb	Sr
Y	Zr	Hf
Nb	Ta	W
Cr	Mn	Fe
Co	Ni	Cu
Zn	Ga	Ge
As	Se	Br
Kr	Rb	Sr
Y	Zr	Nb
Mo	Tc	Ru
Rh	Pd	Ag
Cd	In	Sn
Sb	Te	I
Xe	At	Rn

Drag&Drop your Experimental spectrum and Identify your molecular Formula



# Monoisotopic mass - most abundant isotopes



isotope	%	mass (Da)
$^1\text{H}$	99.99	1.0078
$^{12}\text{C}$	98.93	12.0000
$^{14}\text{N}$	99.64	14.0030
$^{16}\text{O}$	99.76	15.9949

MF	mass (Da)	ppm
NO	29.9979	-2.011
CH <sub>2</sub> O	30.0105	10.565
H <sub>2</sub> N <sub>2</sub>	30.0217	21.798
CH <sub>4</sub> N	30.0343	34.374
C <sub>2</sub> H <sub>6</sub>	30.0469	46.95

# Statistics of pubchem

---



**Pubchem SDF**

Search by  
EM and MF

**Visualizer**

FTP download  
Calculate MF  
MF stats  
Exact mass



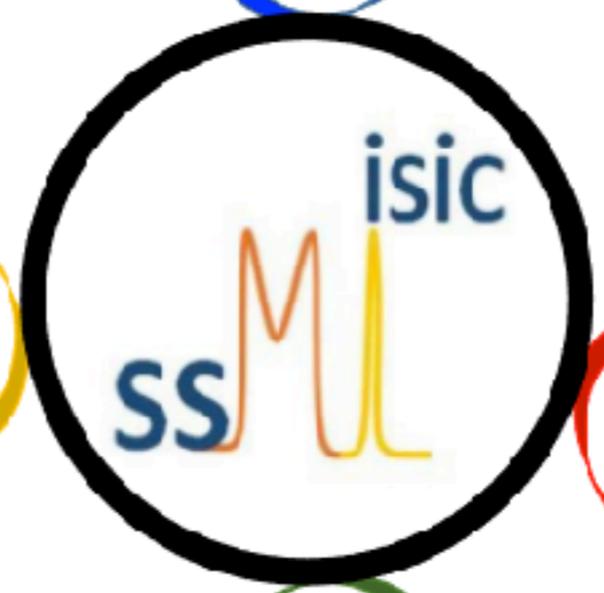
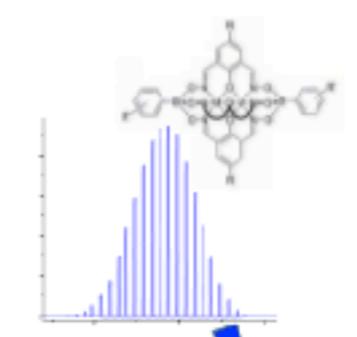
**MongoDB**



**Webservice**

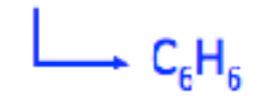


Please cite US

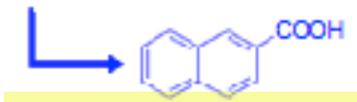


**Calculations**

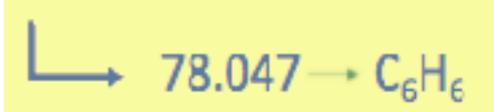
Calculate your theoretical spectrum from a molecular formula



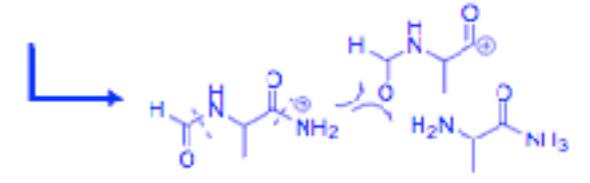
From a chemical structure



Find a MF from a Monoisotopic Mass



MS/MS fragment calculation

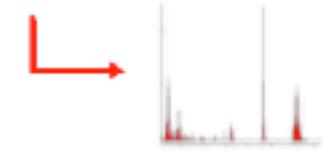


**Similarity to spectrum**

Extend your search using all elements

S	C	N
O	H	F
Br	Cl	I
Na	K	Ca
Fe	Co	Ni
Cu	Zn	Pb

Drag&Drop your Experimental spectrum and Identify your molecular Formula

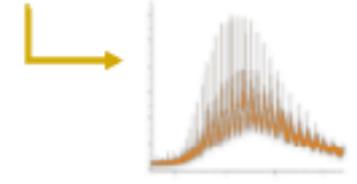


Predict your CID/ETD spectra



**Peptides**

Match your Experimental data



**Match to list**

Drag&Drop and match to contaminants



Try our teaching tools

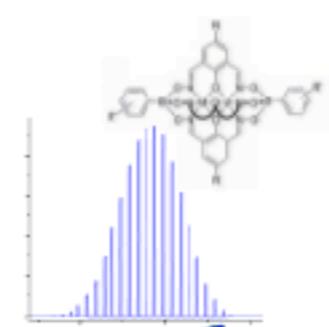
# Contaminants / Custom database

---

- Two possible problems :
  - Mass spectra may contain contaminants
  - You want to look for a list of products
- A solution :
  - Edit a list of molecular formula in a spreadsheet
  - Match directly the list to your spectrum



Please cite US

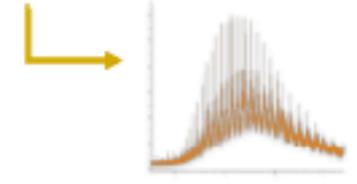


Predict your CID/ETD spectra

Ala-Gly-Ala-Trp

Peptides

Match your Experimental data



Try our teaching tools

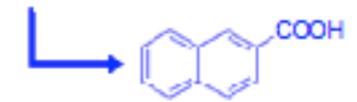
Match to list  
 Drag&Drop and match to contaminants  
 Easycont

Calculations

Calculate your theoretical spectrum from a molecular formula

$C_6H_6$

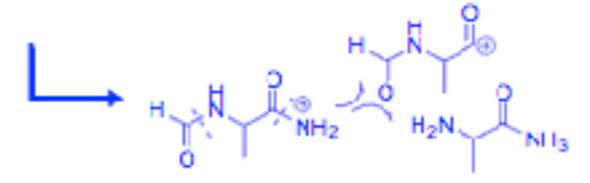
From a chemical structure



Find a MF from a Monoisotopic Mass

78.047  $\rightarrow C_6H_6$

MS/MS fragment calculation

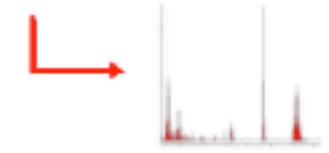


Similarity to spectrum

Extend your search using all elements

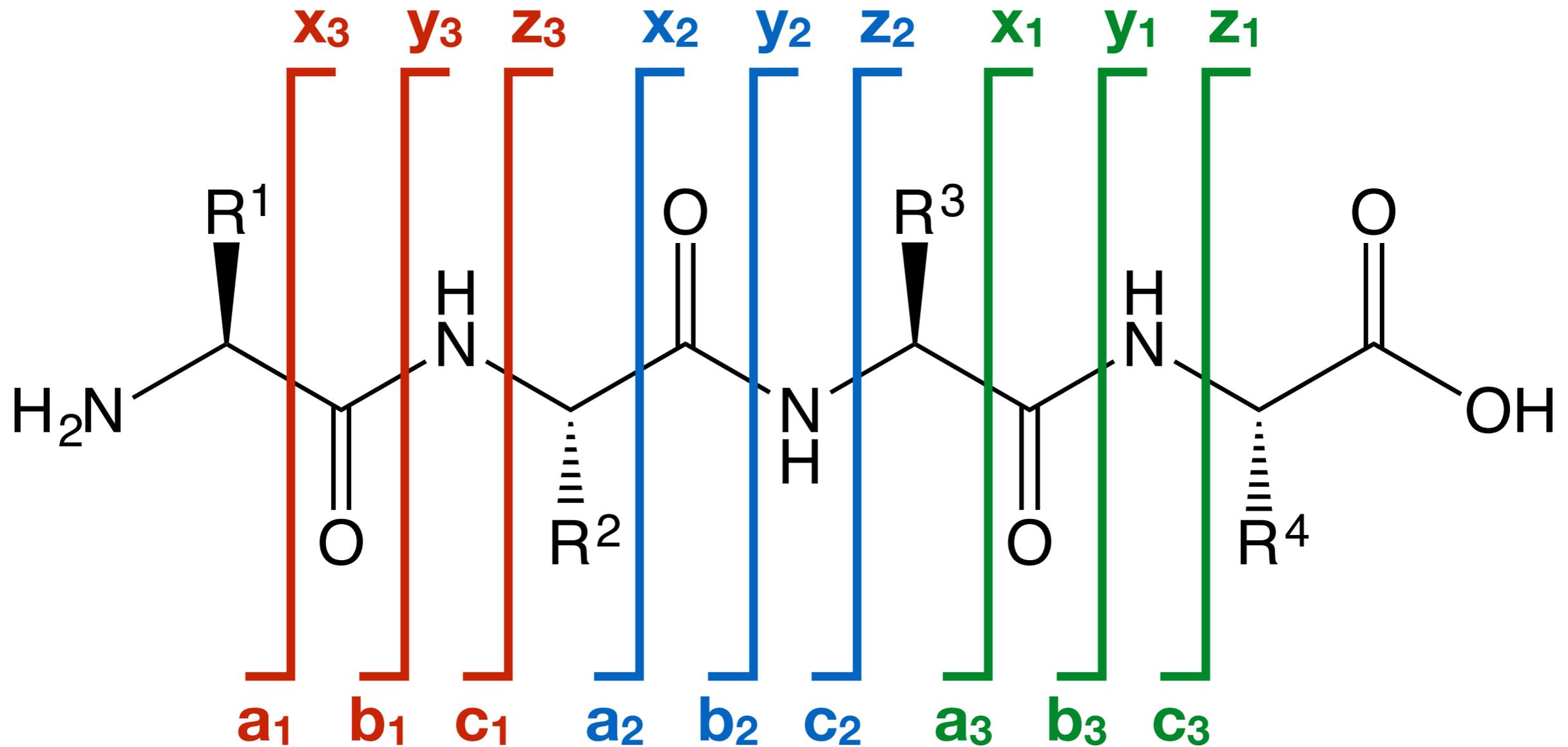
S	C	N
O	H	F
Br	Cl	I
Na	K	Ca
Fe	Pb	As
Se	Ag	Cd
Cr	Mn	Zn
Al	Ga	In
Sn	Bi	Po
At		

Drag&Drop your Experimental spectrum and Identify your molecular Formula



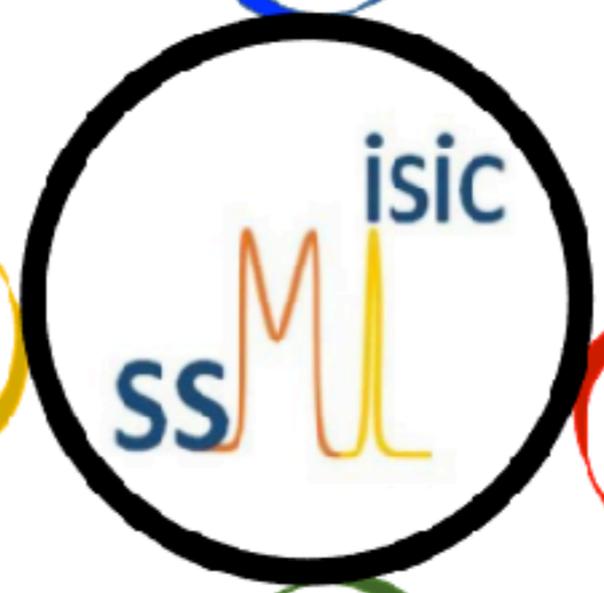
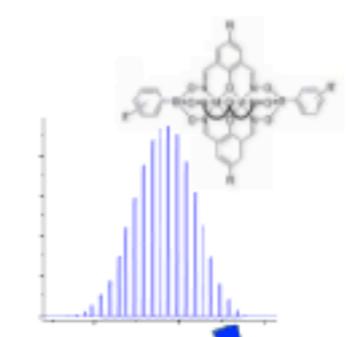
# Peptide fragmentation

---





Please cite US



Predict your CID/ETD spectra  
 ↳ Ala-Gly-Ala-Trp

Match your Experimental data

Try our teaching tools



Match to list  
 Drag&Drop and match to contaminants  
 ↳ Easycont

Calculations

Calculate your theoretical spectrum from a molecular formula  
 ↳  $C_6H_6$

From a chemical structure  
 ↳

Find a MF from a Monoisotopic Mass  
 ↳ 78.047 →  $C_6H_6$

MS/MS fragment calculation  
 ↳

Similarity to spectrum

Extend your search using all elements  
 ↳

Drag&Drop your Experimental spectrum and Identify your molecular Formula  
 ↳

# Cheminformatics

---



# From 2D drawing to conformations

---

## Draw

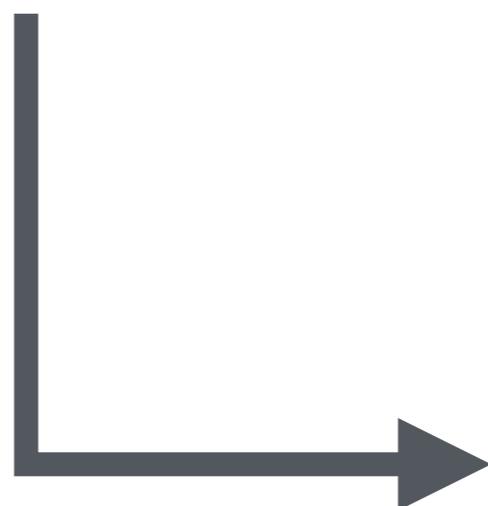
2D molfile

## Display

conformers

## Server

moloc  
2D to 3D

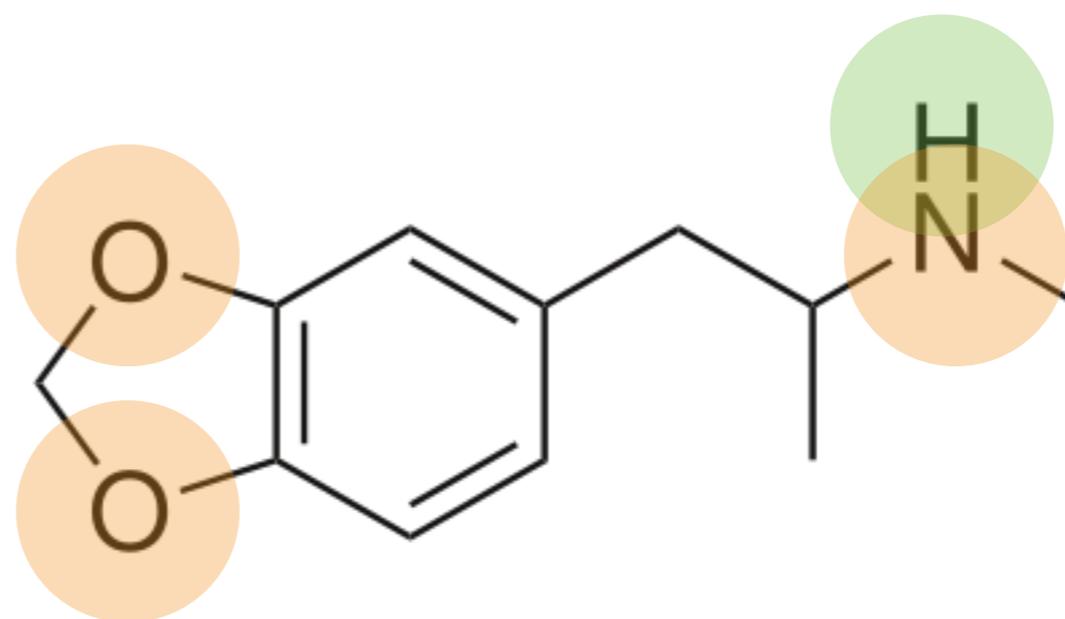


Demo

# Lipinski rule of 5

---

- ✓ H bond donors  $\leq 5$
- ✓ H bond acceptors  $\leq 10$
- ✓  $\log P \leq 5$  **1.8**
- ✓  $mw \leq 500$  **193**



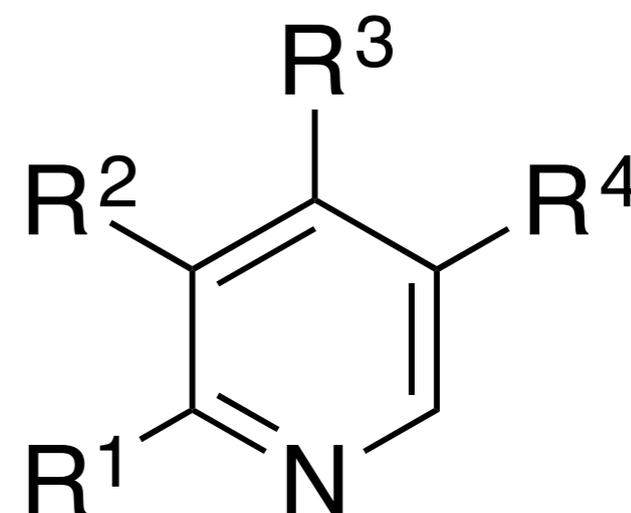
Demo 1

Demo 2

# Virtual combinatorial library

---

- The problem
  - Search for products with specific properties :
    - logP
    - mw
    - H donor
    - H acceptor
- A solution
  - Generate all possible molecules based on a template and fragments
  - Predict the properties
  - Filter using parallel coordinates



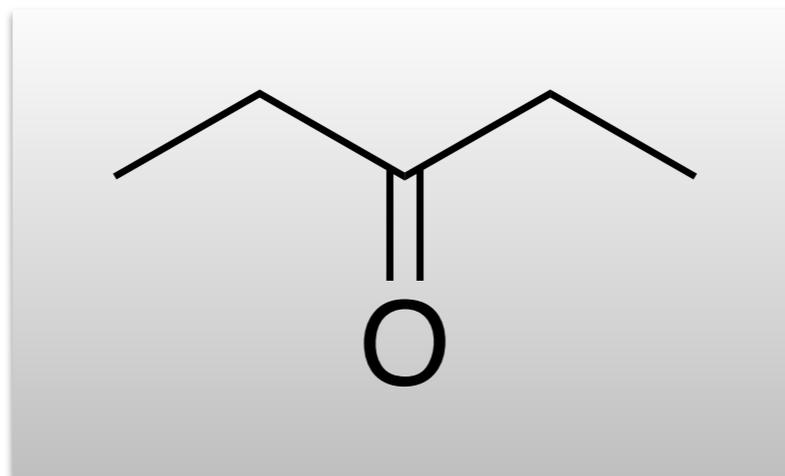
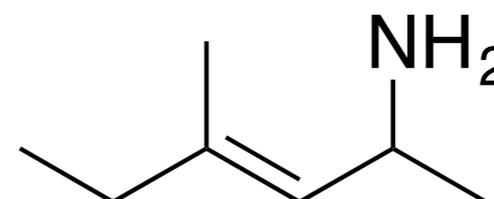
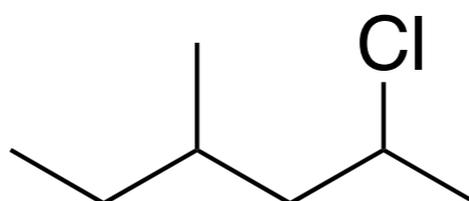
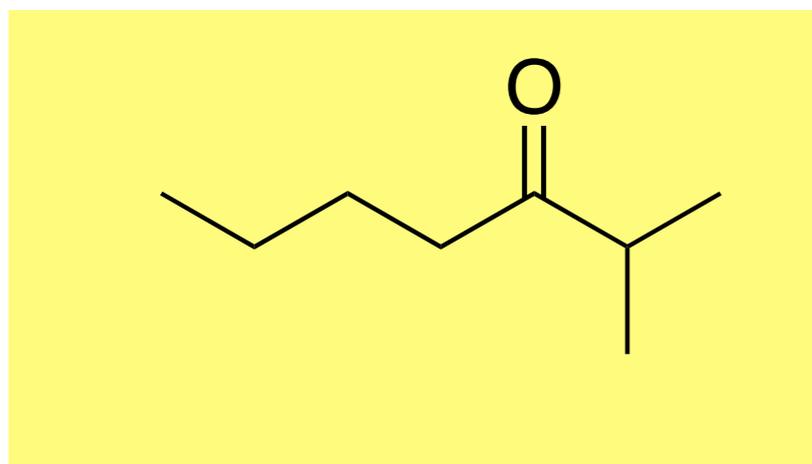
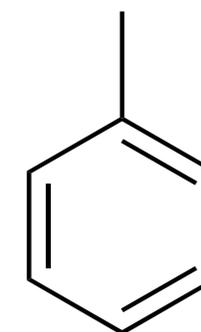
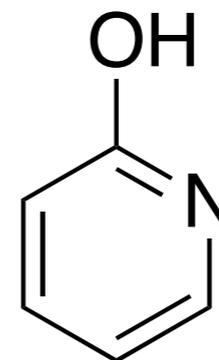
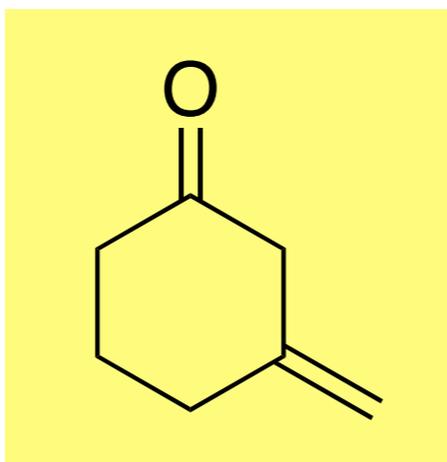
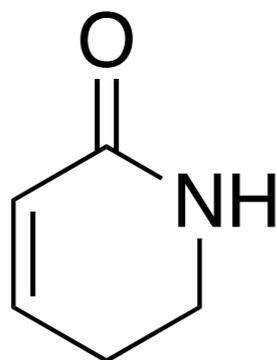
Demo

# Substructure search - Wikipedia

---

# Substructure search ?

---



W Benzene - Wikipedia, the free encyclopedia Luc Faliny

← → ↻ ↗ <https://en.wikipedia.org/wiki/Benzene> ★ ABP 🔍 ☰

Create account Log in



**WIKIPEDIA**  
The Free Encyclopedia

- [Main page](#)
- [Contents](#)
- [Featured content](#)
- [Current events](#)
- [Random article](#)
- [Donate to Wikipedia](#)
- [Wikipedia store](#)

Interaction

- [Help](#)
- [About Wikipedia](#)
- [Community portal](#)
- [Recent changes](#)
- [Contact page](#)

Tools

- [What links here](#)
- [Related changes](#)
- [Upload file](#)
- [Special pages](#)
- [Permanent link](#)
- [Page information](#)
- [Wikidata item](#)
- [Cite this page](#)

Print/export

- [Create a book](#)
- [Download as PDF](#)
- [Printable version](#)

Article Talk

## Ber

From Wikipedia, the free encyclopedia

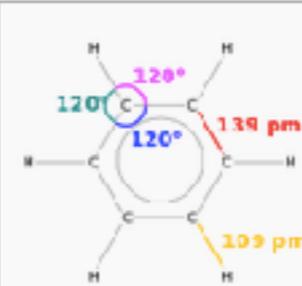
**Benzene** (<sup>*Talk*</sup> <sup>*Not on Wikidata*</sup>)

Names	
IUPAC name	Benzene
Systematic IUPAC name	Cyclohexa-1,3,5-triene
Other names	1,3,5-Cyclohexatriene, Benzol, Phene, Phenyl hydride
Identifiers	
CAS Registry Number	71-43-2 ✓
ChEBI	CHEBI:16716 ✓
ChEMBL	ChEMBL277500 ✓
ChemSpider	236 ✓
EC number	200-753-7
InChI	<span>[show]</span>
Jmol-3D images	<span>Image</span> <span><span><span></span></span></span>
KEGG	C01407 ✓
PubChem	241
RTECS number	CY1400000
SMILES	<span>[hide]</span>
1 History	<span><b>c1ccccc1</b></span>
1 InChI	J64922108F ✓
Properties	
1.3 Early applications	
2 Structure	
3 Benzene derivatives	
4 Production	

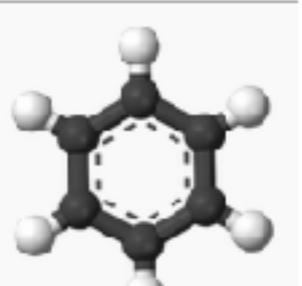
*s, see Benzene (disambiguation).*

chemical formula  
ring, with 1 hydrogen  
contain only carbon

most elementary  
second *[n]*-annulene  
. It is sometimes  
be liquid with a sweet  
h as ethylbenzene  
because it has a high  
rising a few percent  
by benzene's



Geometry of molecule



Ball and stick model of molecule



Benzene molecule

Names
IUPAC name
Benzene
Systematic IUPAC name
Cyclohexa-1,3,5-triene

Read Edit View history

Search

c1ccccc1

Every night ...

---

## Find

ChemBox  
DrugBox



## Extract

SMILES



## Generate

JSON



## Publish

GitHub

<https://github.com/cheminfo/wikipedia>  
<http://www.cheminfo.org/wikipedia>

Demo

# Simple database of chemicals - open source malaria

---

## Store

Google  
spreadsheet



## Display

visualizer

<http://malaria.cheminfo.org>

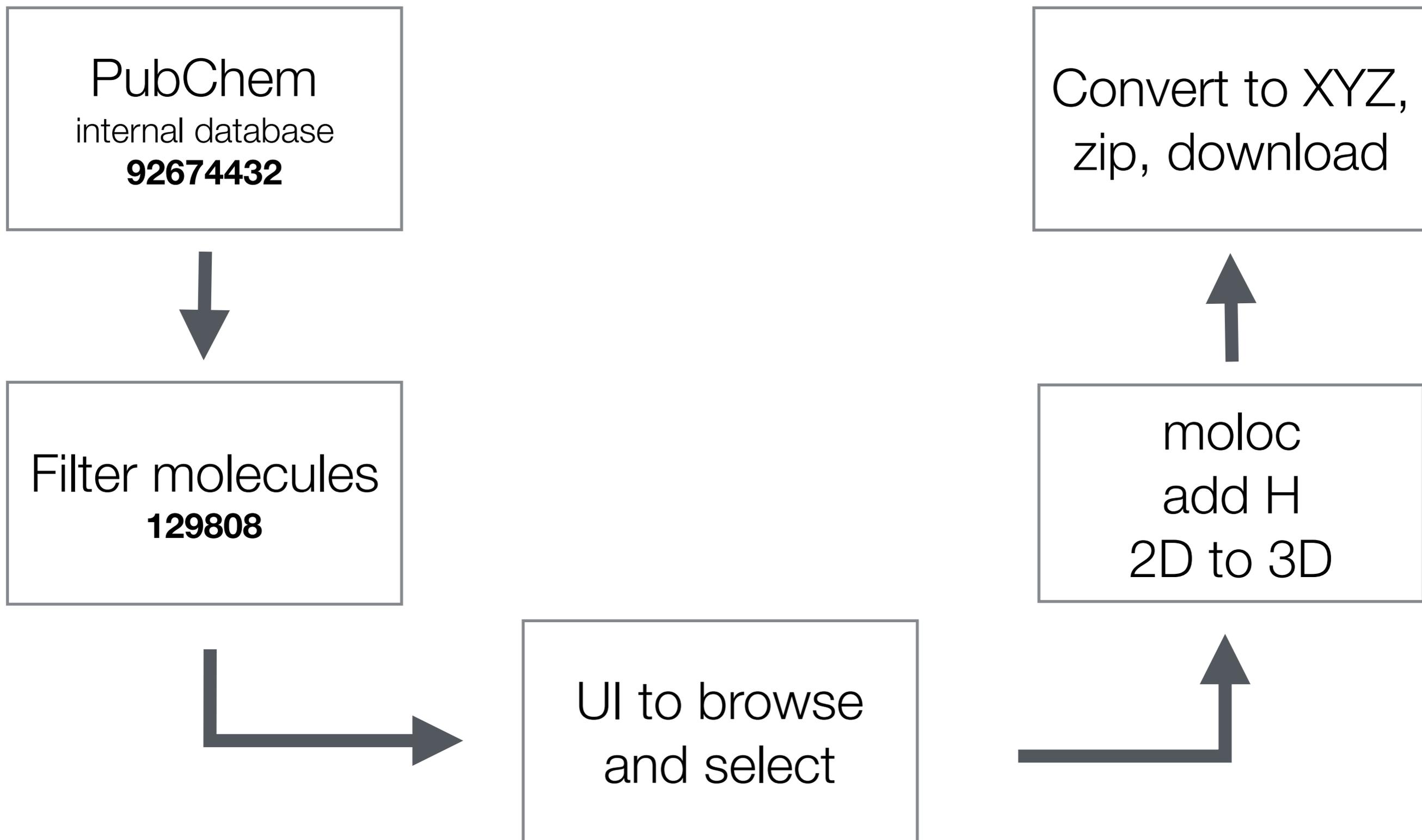
"Open source drug discovery: highly potent antimalarial compounds derived from the tres cantos arylpyrroles."  
*ACS central science* 2, no. 10 (2016): 687-701.

[Link to google docs](#)

[To the visualizer](#)

# Validation dataset for ab initio project

---



Demo

# Image analysis

---

# Image analysis : Bacteria in water

---

Coliscan  
Easygel

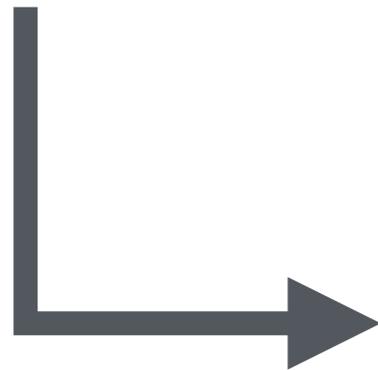


Image  
analysis

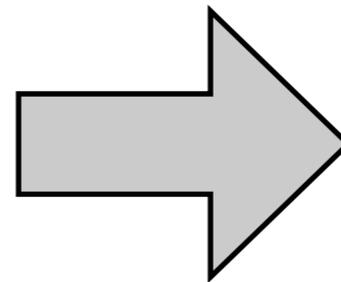
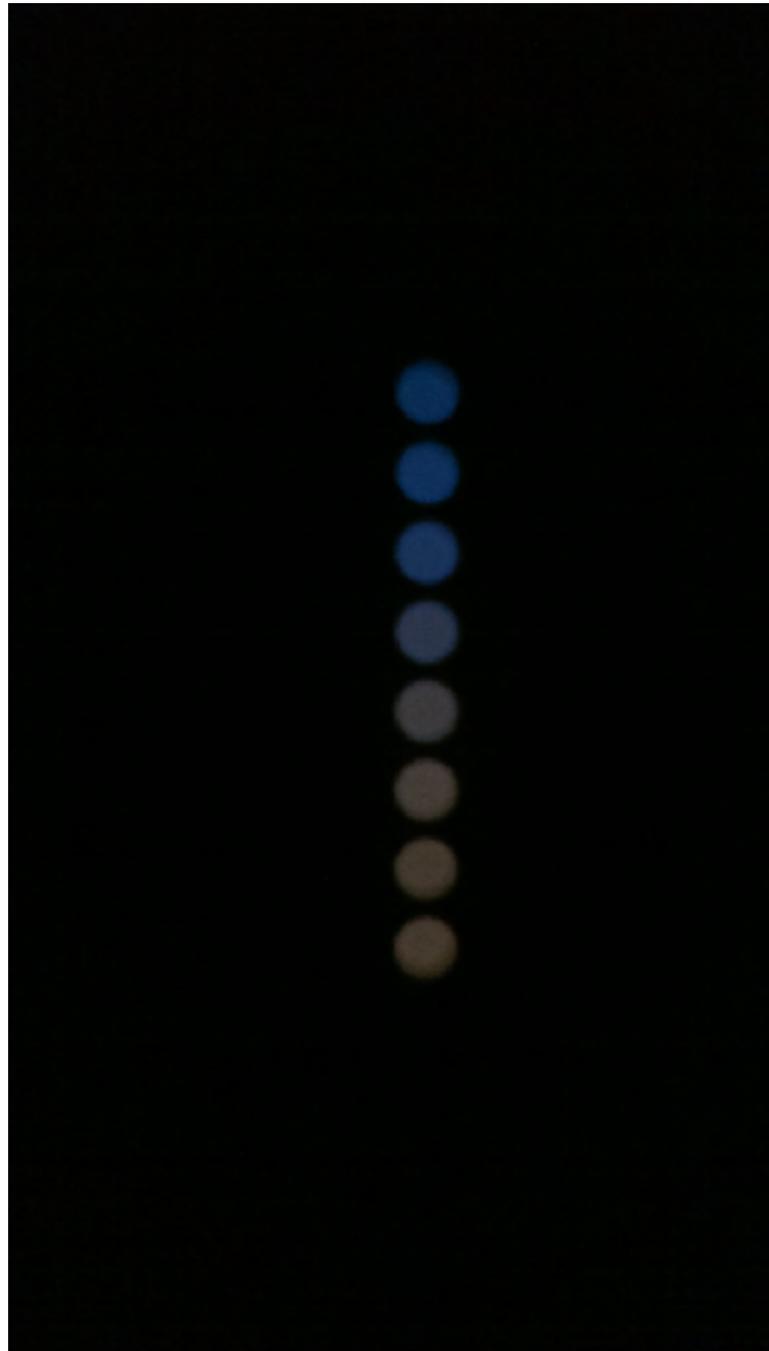
Classification



Demo

# Image analysis: IC<sub>50</sub>

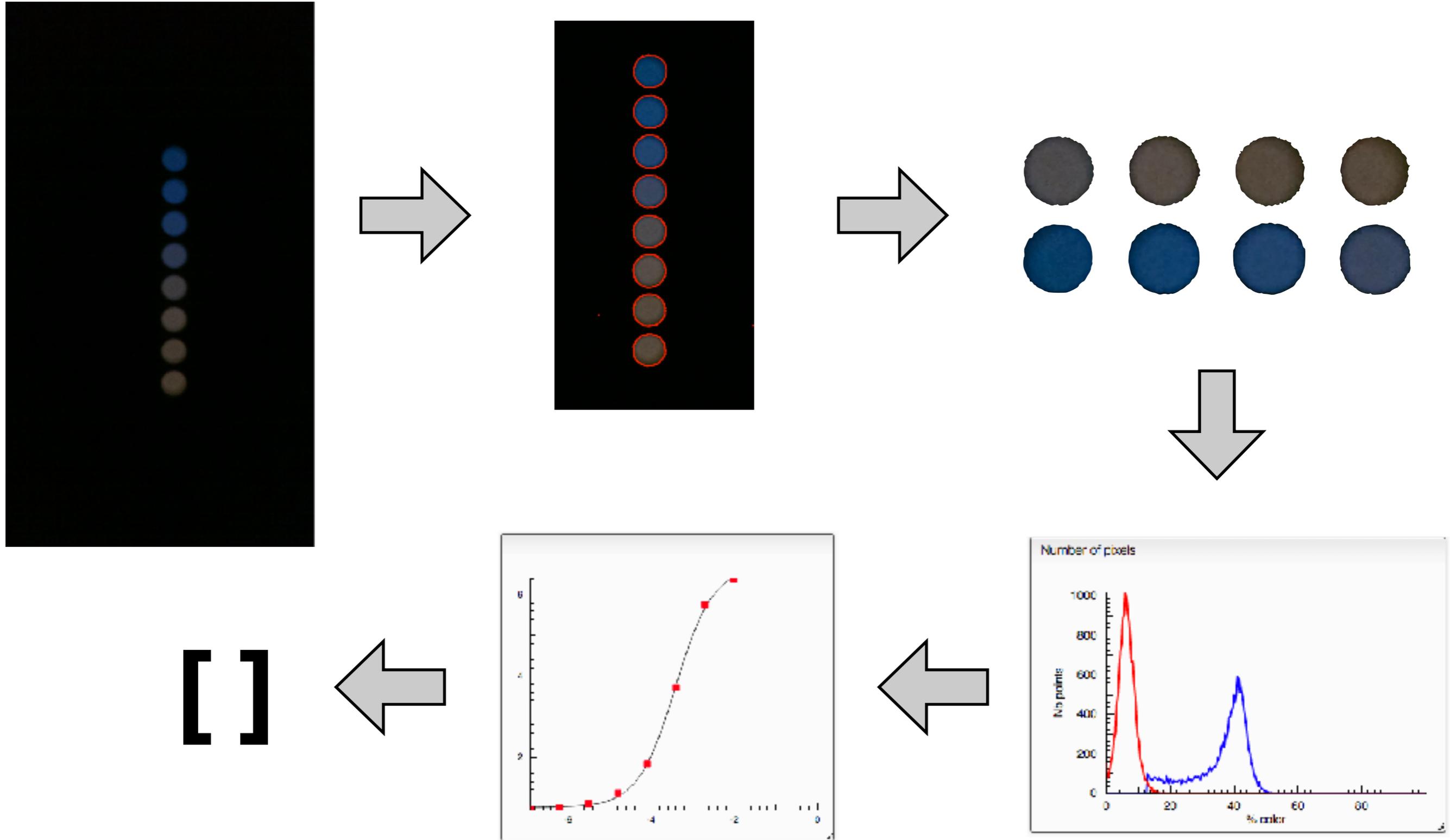
---



[ ]

Bioluminescent sensor proteins for point-of-care therapeutic drug monitoring  
Kai Johnsson *et al.*, Nature chemical biology 2014, 10, 598-603.

# The problem : many steps



Demo

# SEM / TEM image analysis

---

TEM image

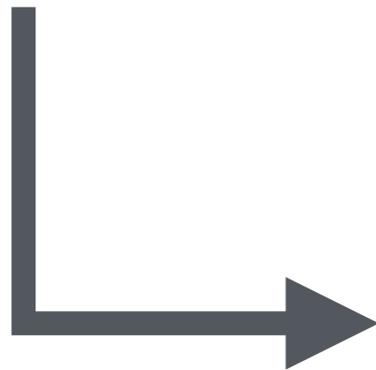


Image  
analysis



Statistics

Demo

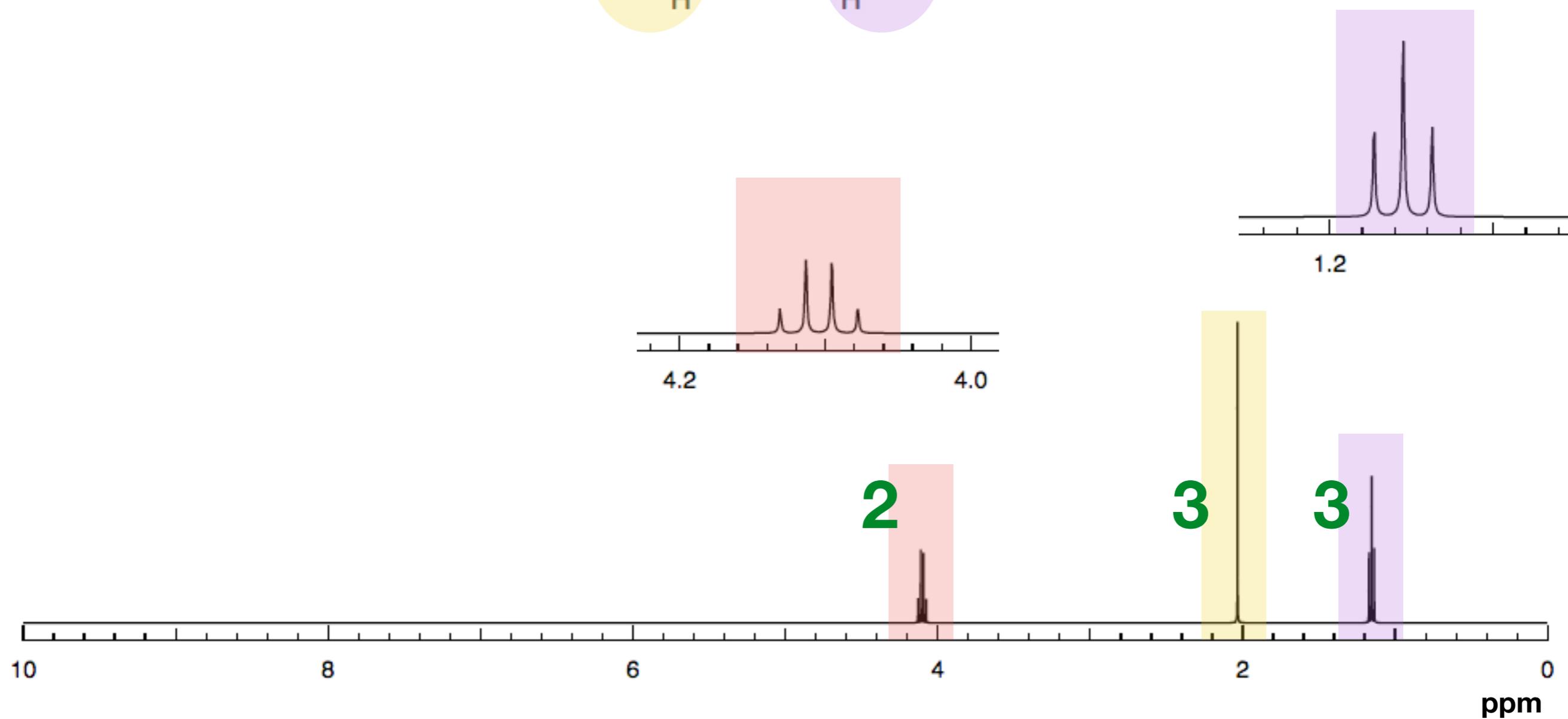
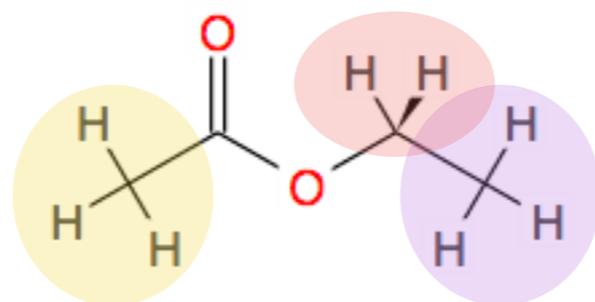
# **Nuclear magnetic resonance**

---

<http://www.nmrdb.org>

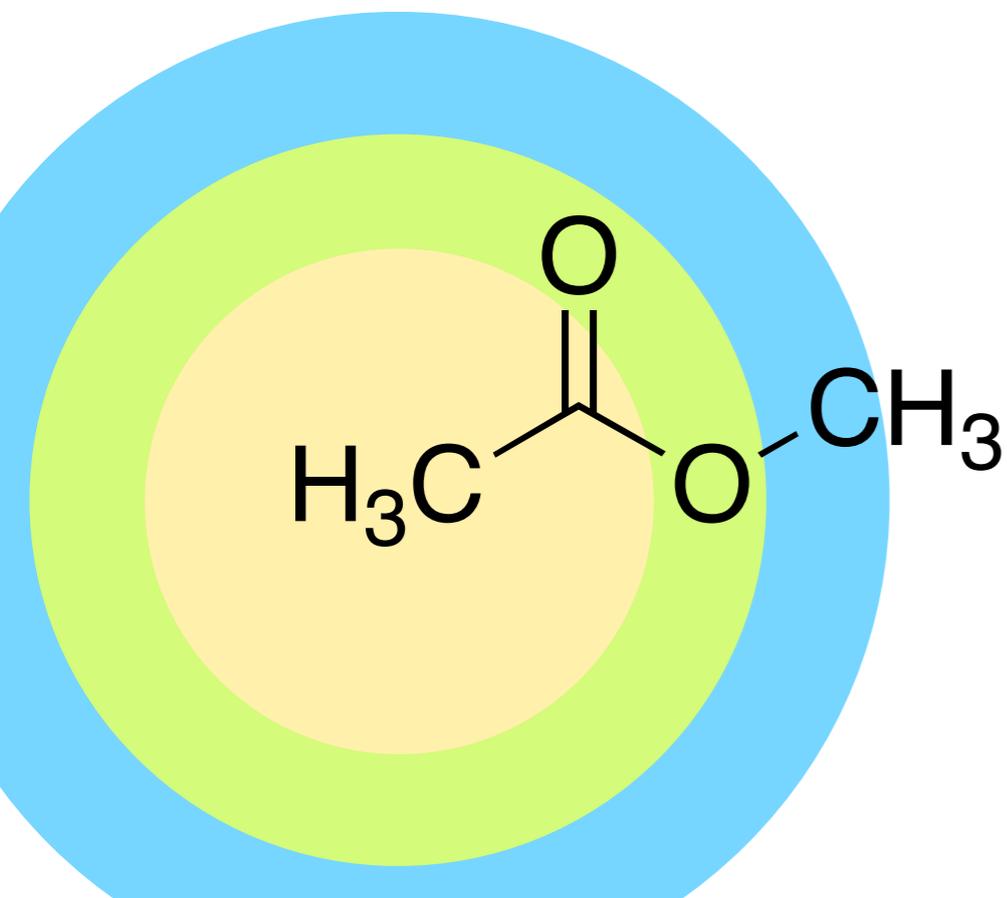
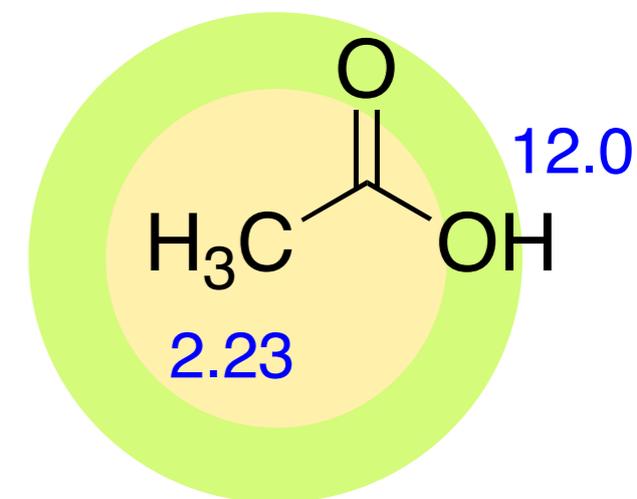
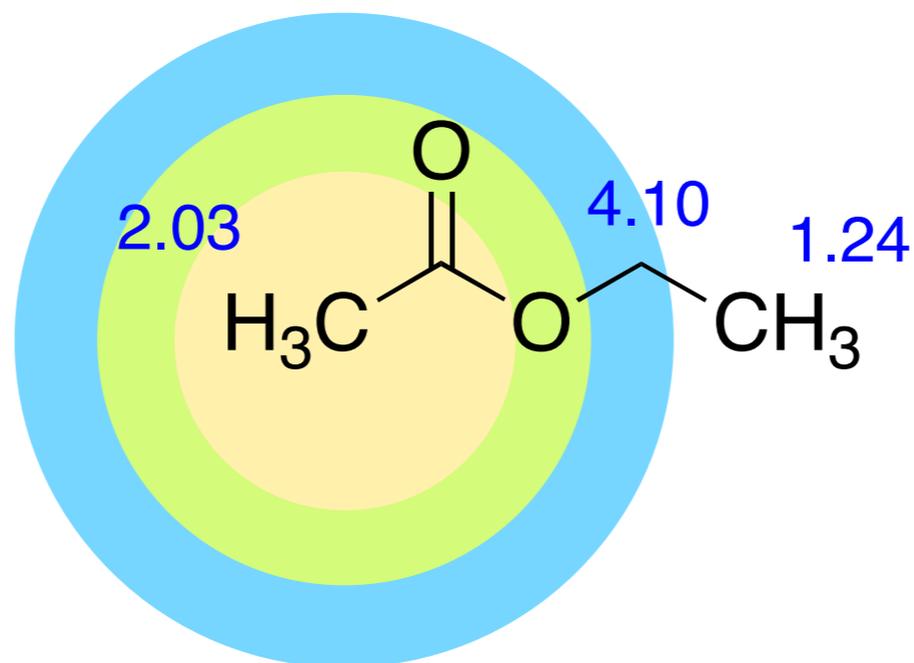
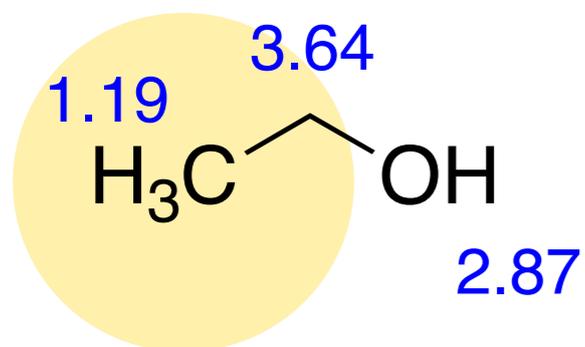


# $^1\text{H}$ NMR spectroscopy



# HOSE code

---



$$\delta = \frac{1.19 + 2.03 + 2.23}{3} = 1.81$$

$$\delta = \frac{2.03 + 2.23}{2} = 2.13$$

$$\delta = 2.03$$

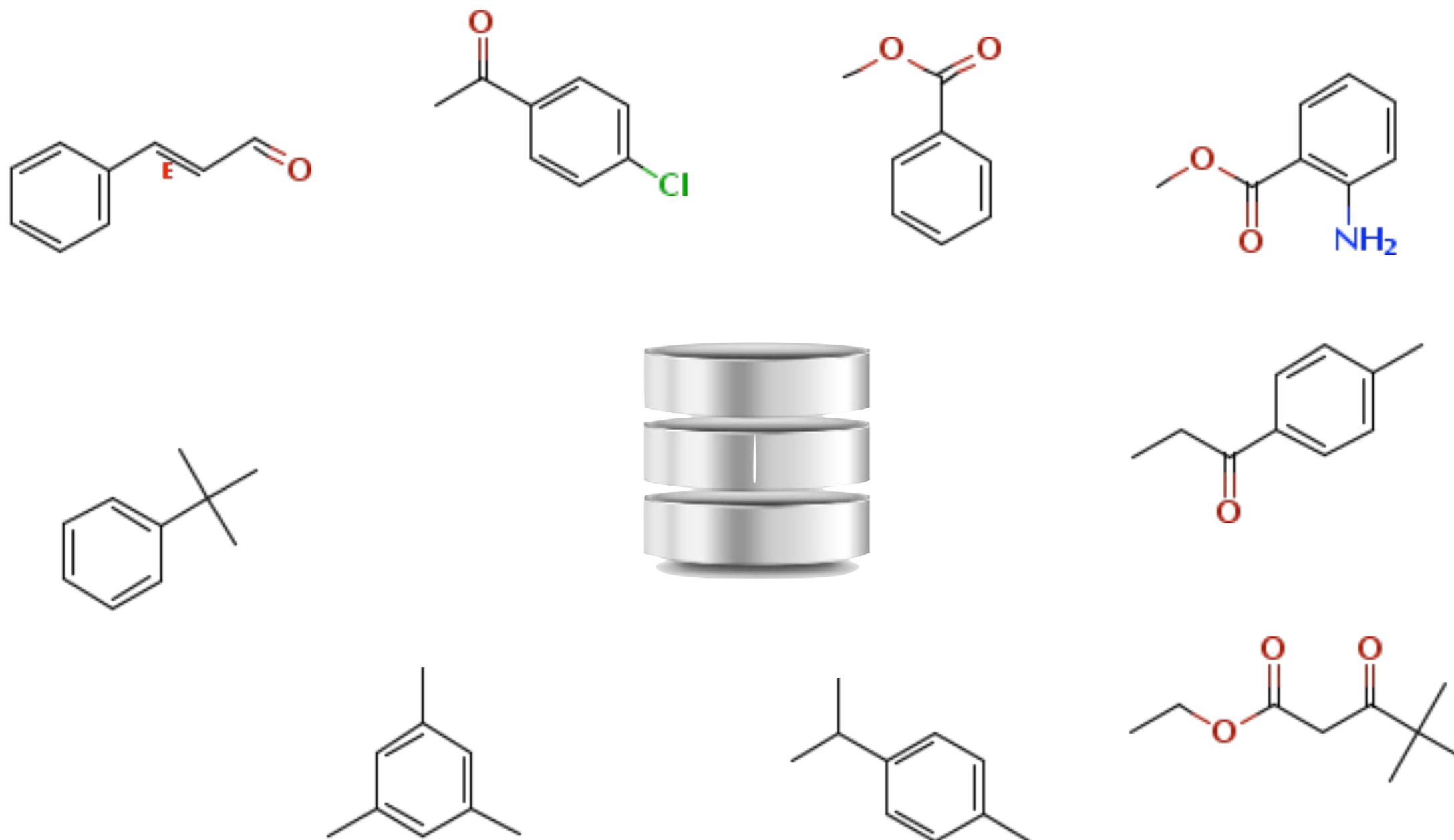
# NMR self learning algorithm - Ask Ernö

---

Creating a NMR chemical shifts predictor without using chemical shifts !?

# 1. Creating a dataset: molfile / spectrum

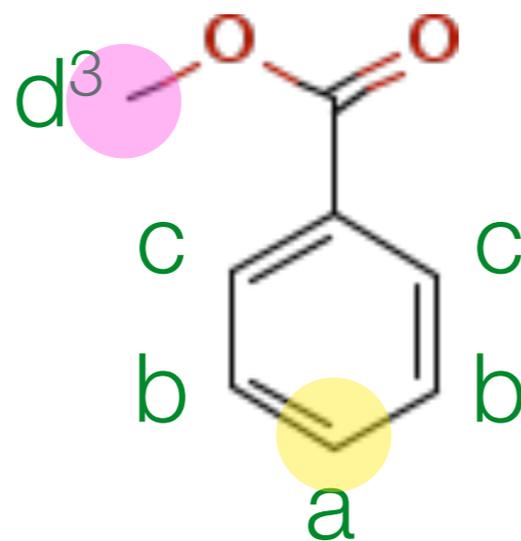
---



2341 molecules / NMR spectra

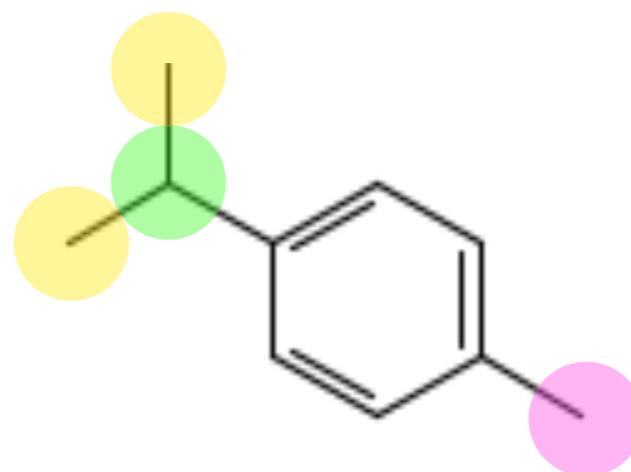
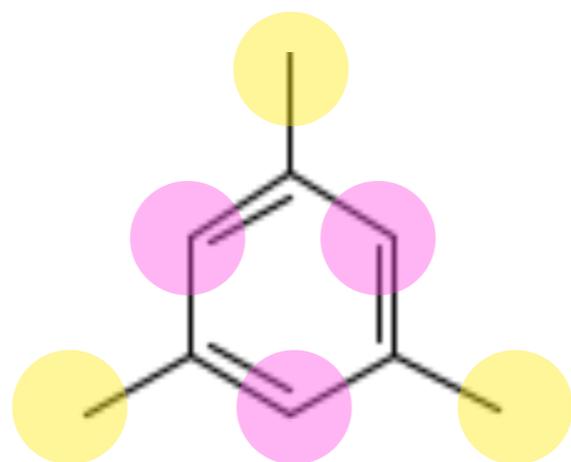
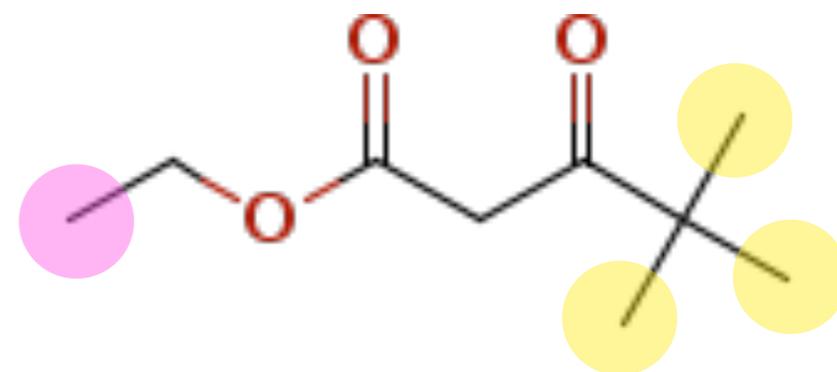
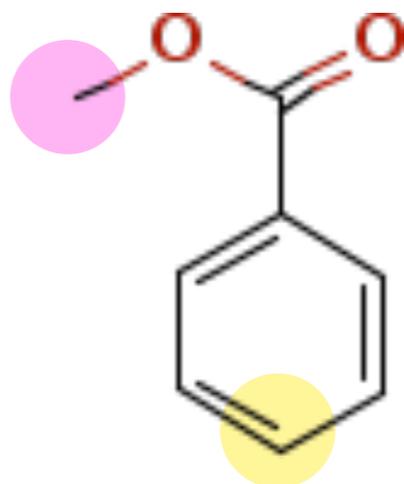
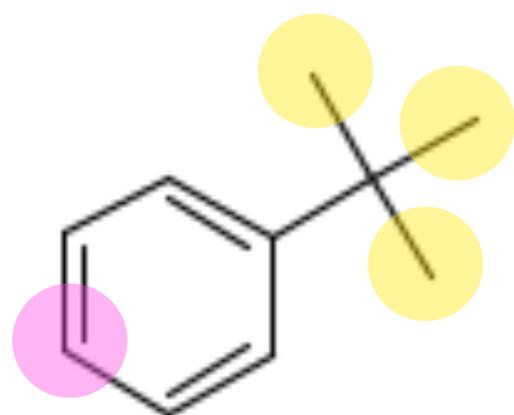
## 2. Number and kind of hydrogens

---



### 3. Unambiguous number of protons

---



### 3. Peak picking and automatic integration

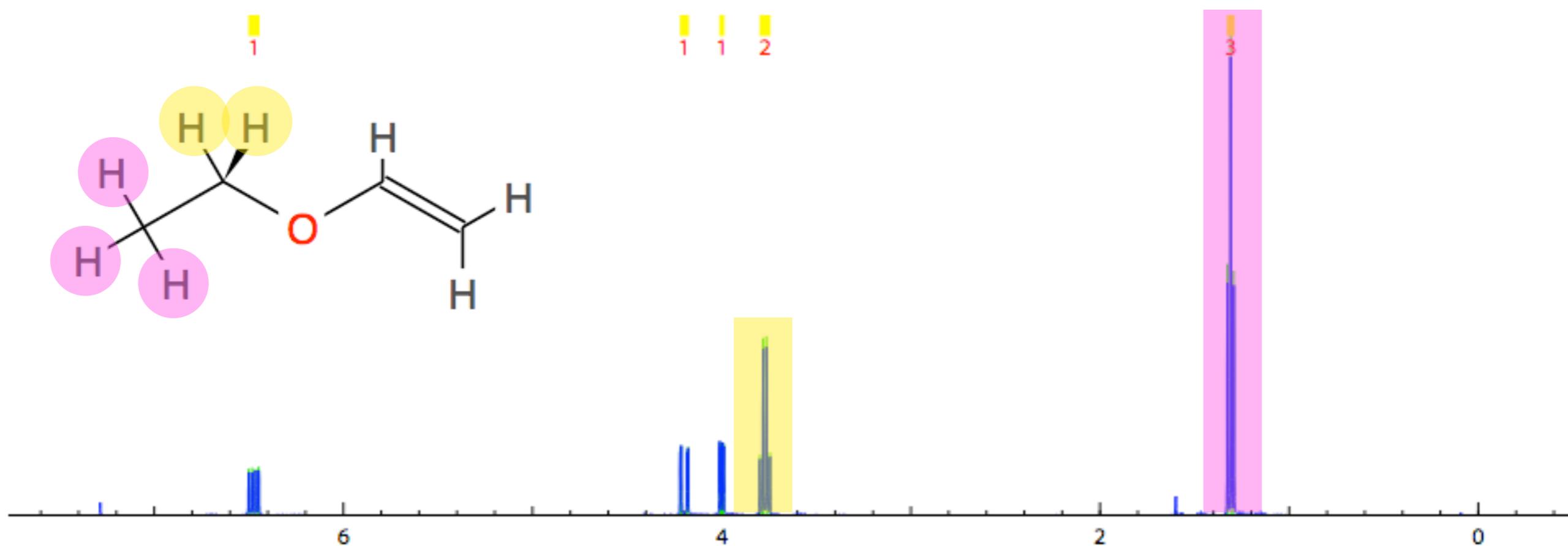
---

- Integration of the zones
- Removal of the NMR solvent / impurities

Demo

## 4. Assign non ambiguous protons

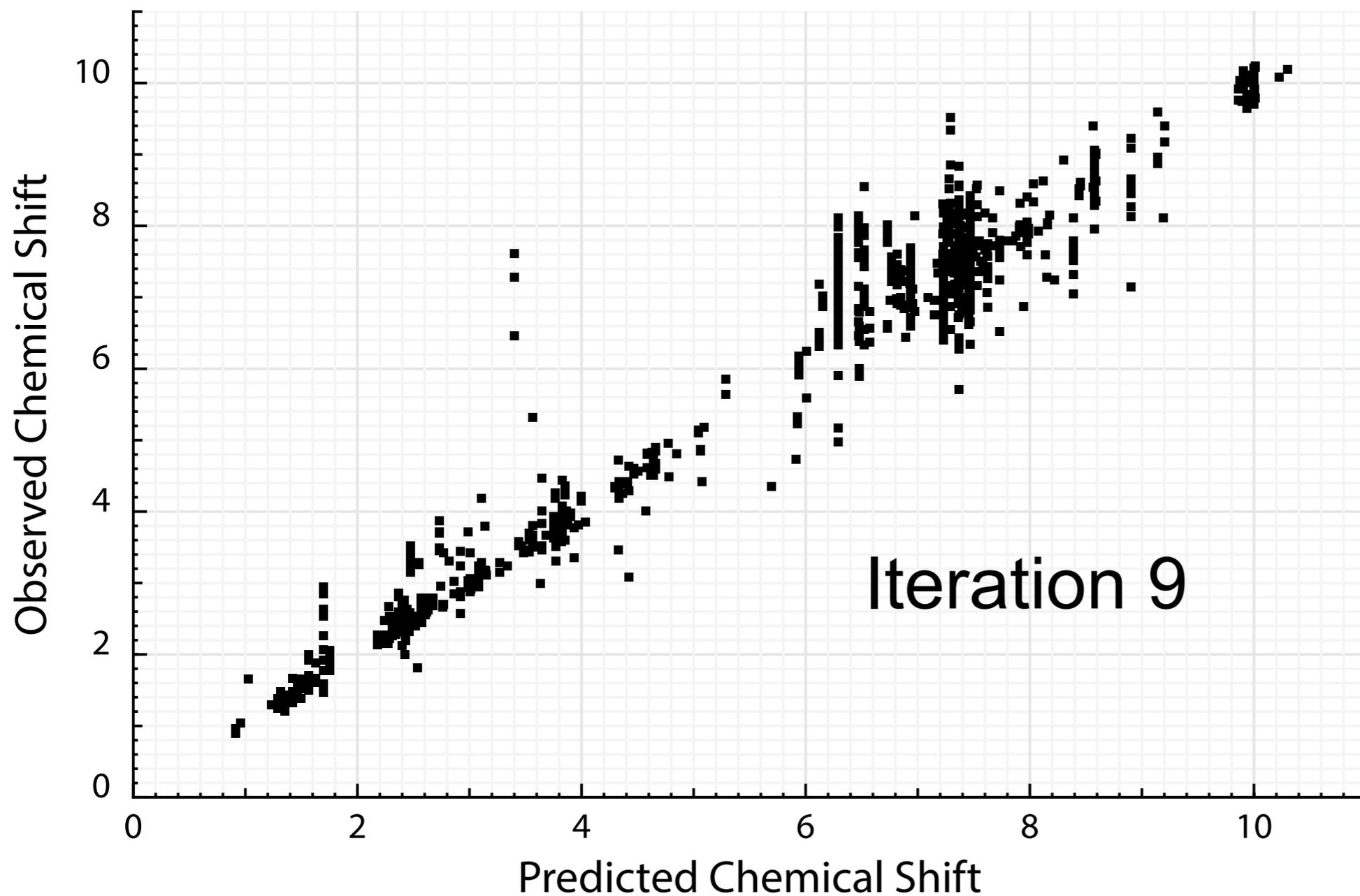
---



... and create corresponding HOSE codes

# Analysis of the 2341 molecule / spectra set

---

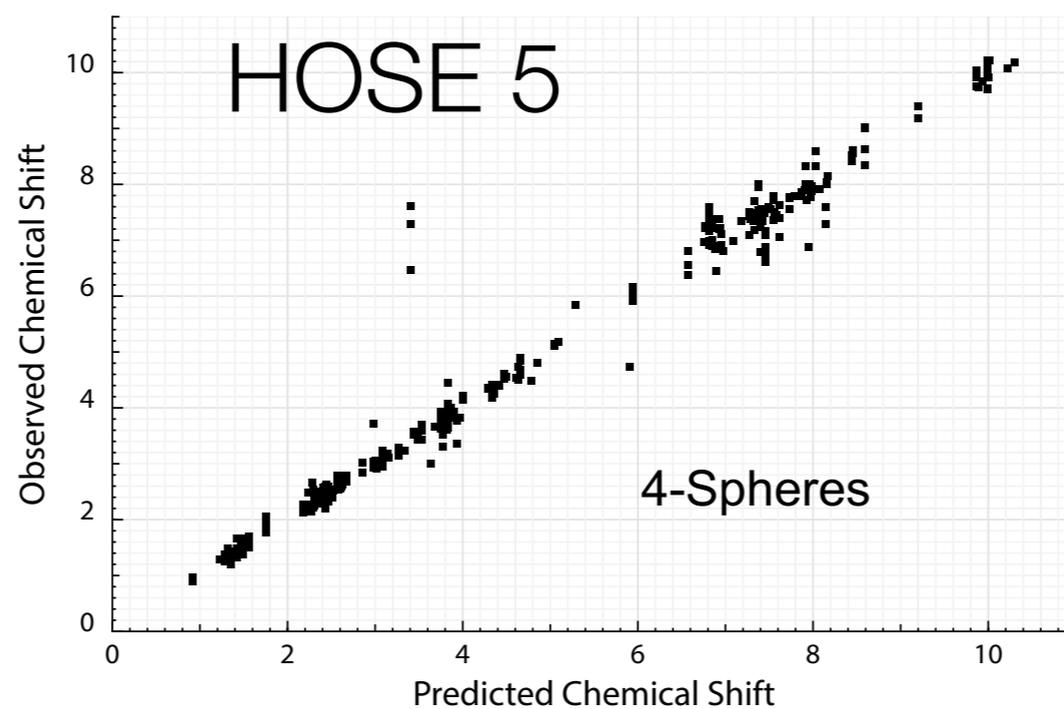
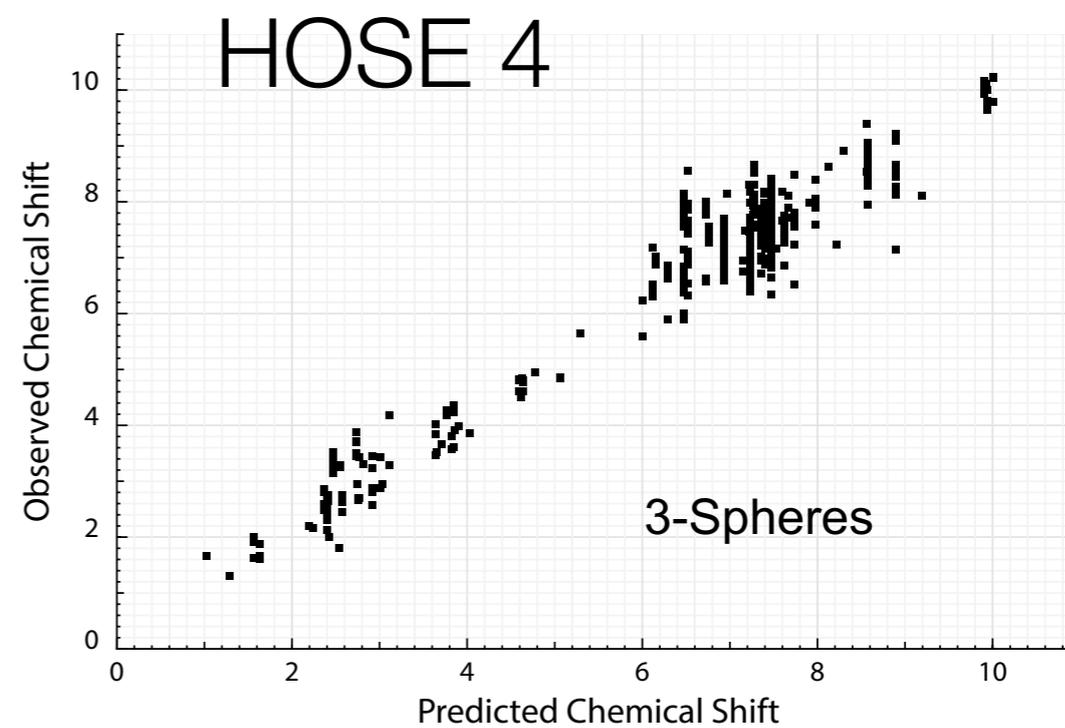
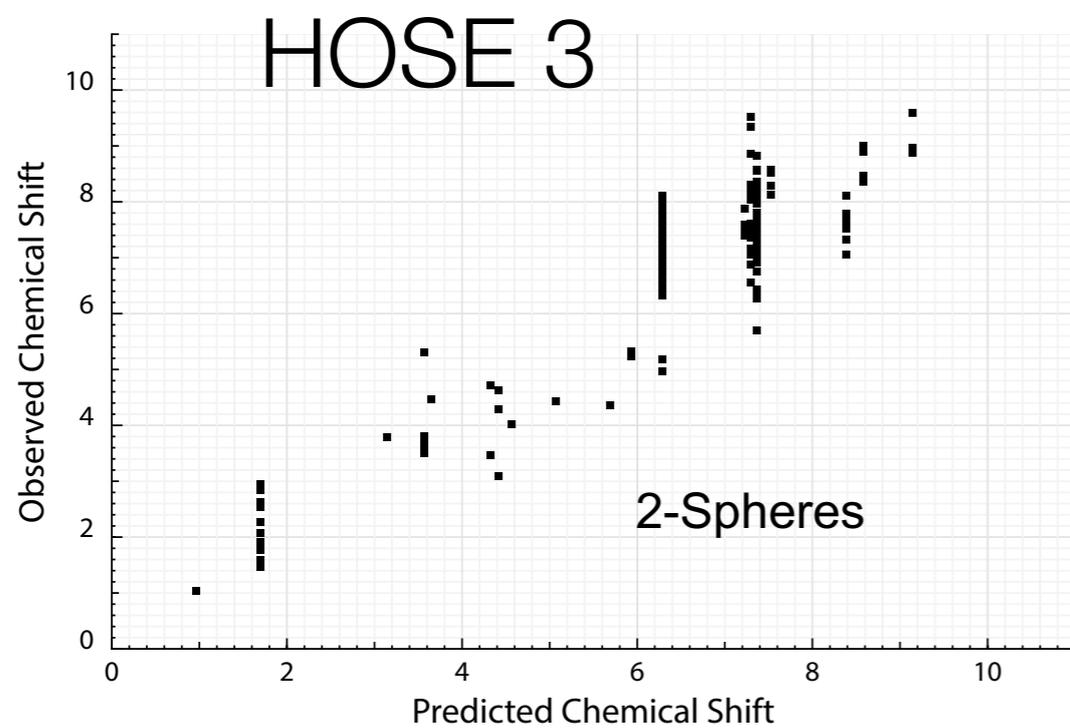


HOSE 3 : 63  
HOSE 4 : 382  
HOSE 5 : 916

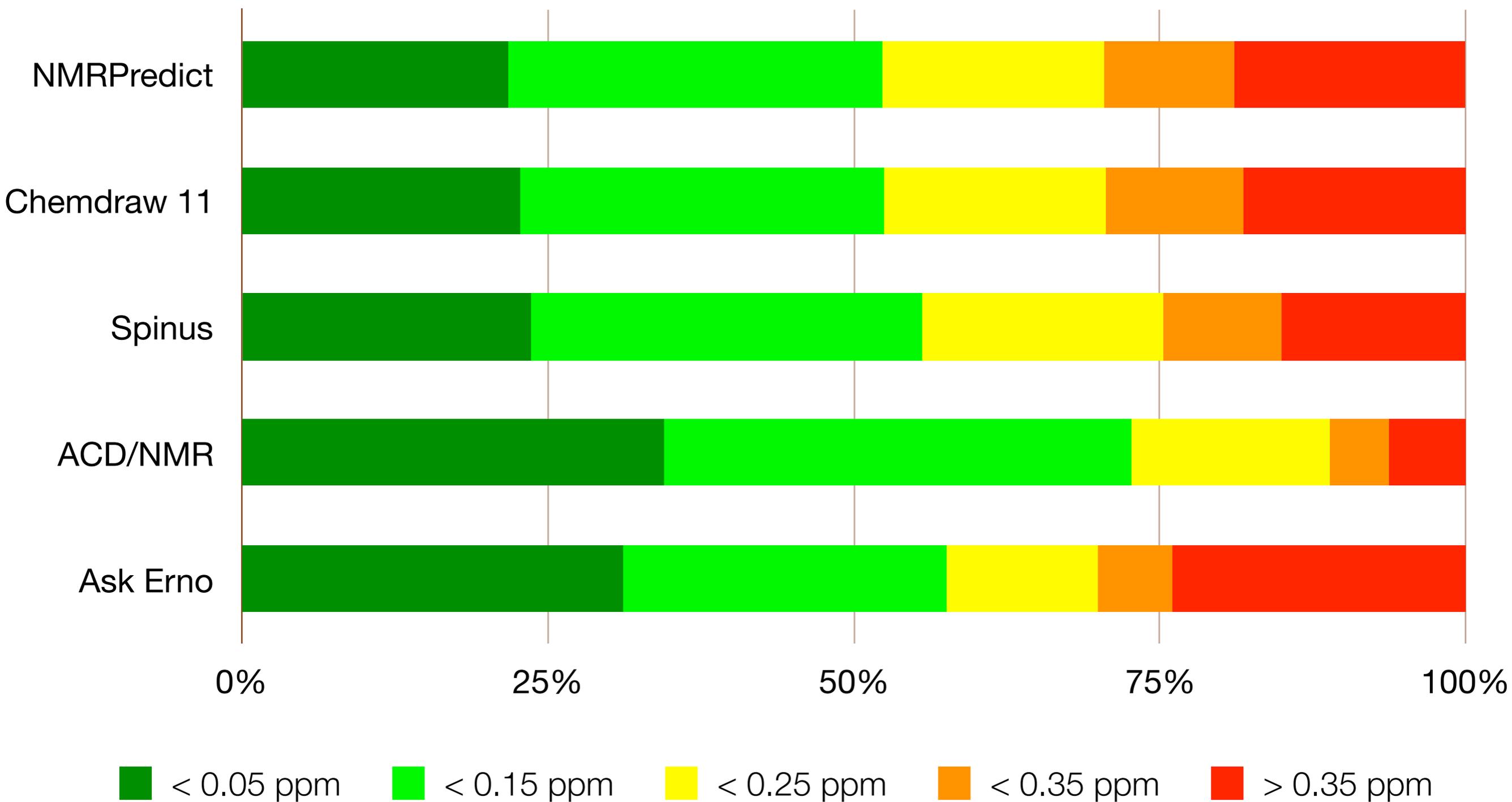
Test set of 298 molecules

# Quality of prediction based on HOSE code level

---



# Prediction average errors

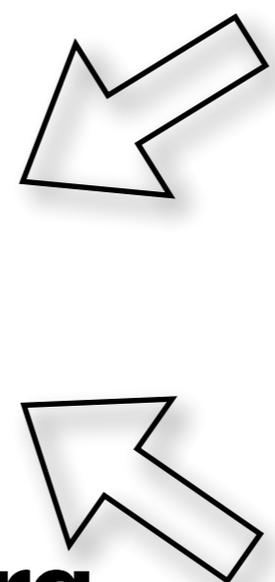
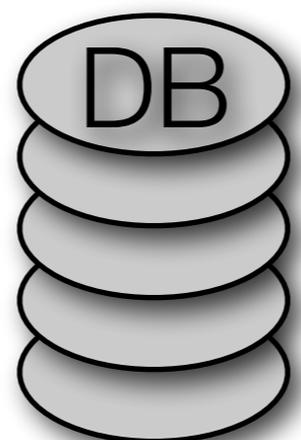
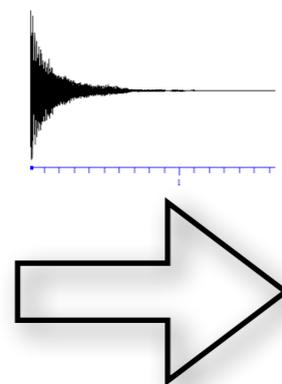




**We need reliable structured data !**

---

<http://www.c6h6.org>



[www.c6h6.org](http://www.c6h6.org)

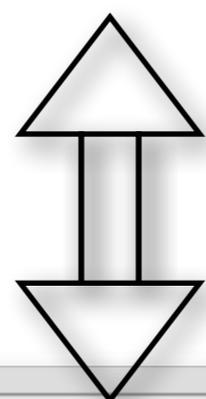
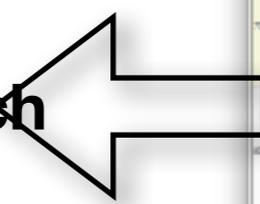
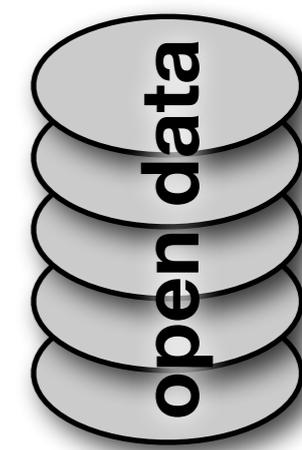
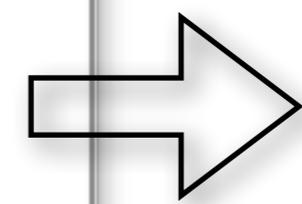


Image analysis  
Mass analysis  
Similarity search  
Predictions

...



ID	Date	Chemical Structure
1234	2013-01-15	<chem>CC1=CC=CC=C1</chem>
5678	2013-01-20	<chem>CC(=O)OC</chem>
9012	2013-02-01	<chem>CC1=CC=C(C=C1)O</chem>



Demo

# Conclusions

---

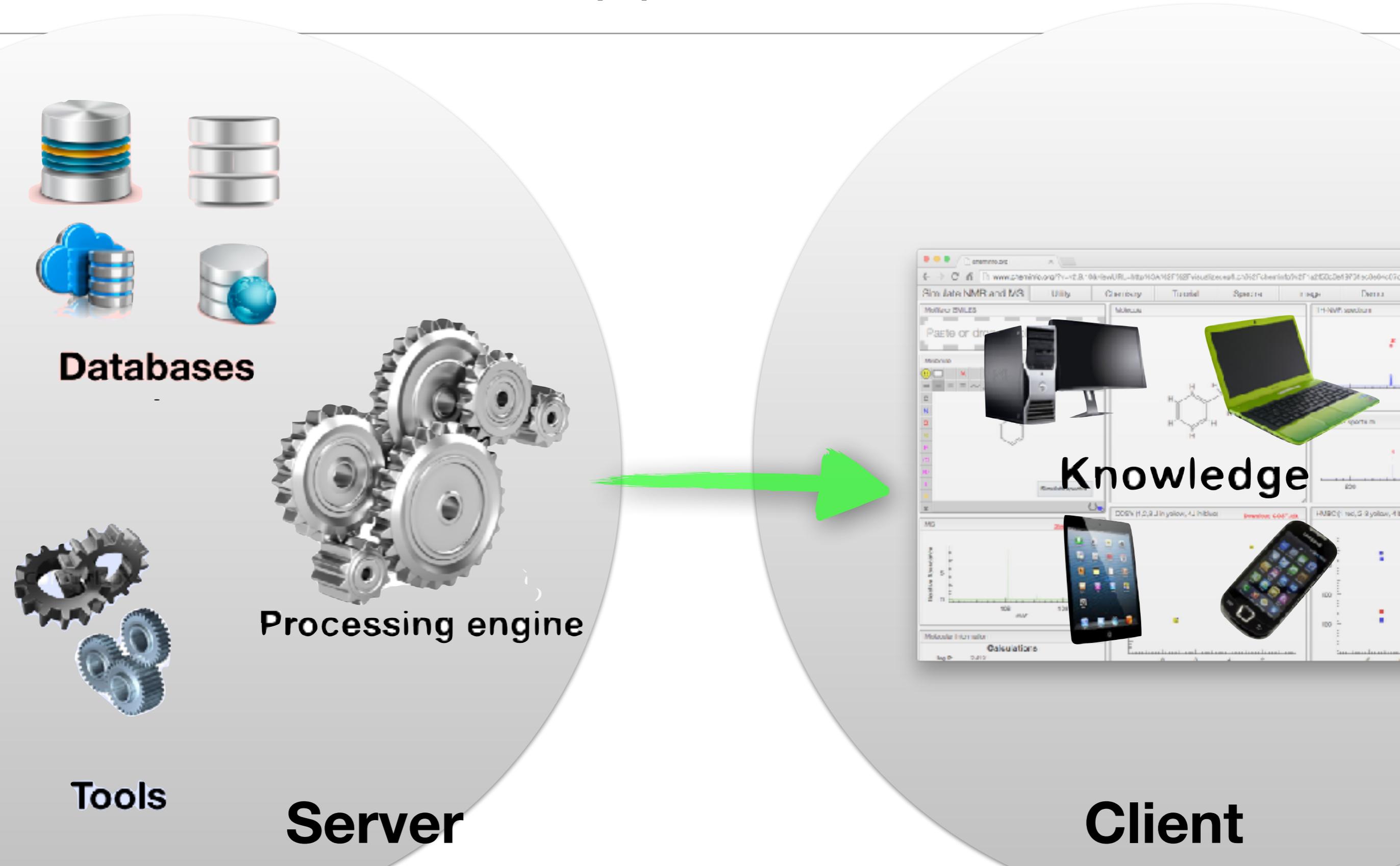


# Advanced visualization platform

---

- Solves scientific problems in the browser
- Pure HTML5 / javascript
- Easily reusable javascript library
- Research, service and teaching
- Over 150 tools in 3 years

# Client / Server approach



# Thanks !

---

## ChemCalc

- Alain Borel (Library EPFL)
- Laure Menin
- Michael Krompiec
- Marek Noga

## NMRdb

- Julien Wist
- Andrès Castillo
- Andrès Bernal

## script, visualizer, components

- Norman Pellet
- Michaël Zasso
- Daniel Kostro
- Jefferson Hernández
- Miguel Asencio
- Julien Wist
- Andrès Castillo



