



Finding Small Molecules In Big Data

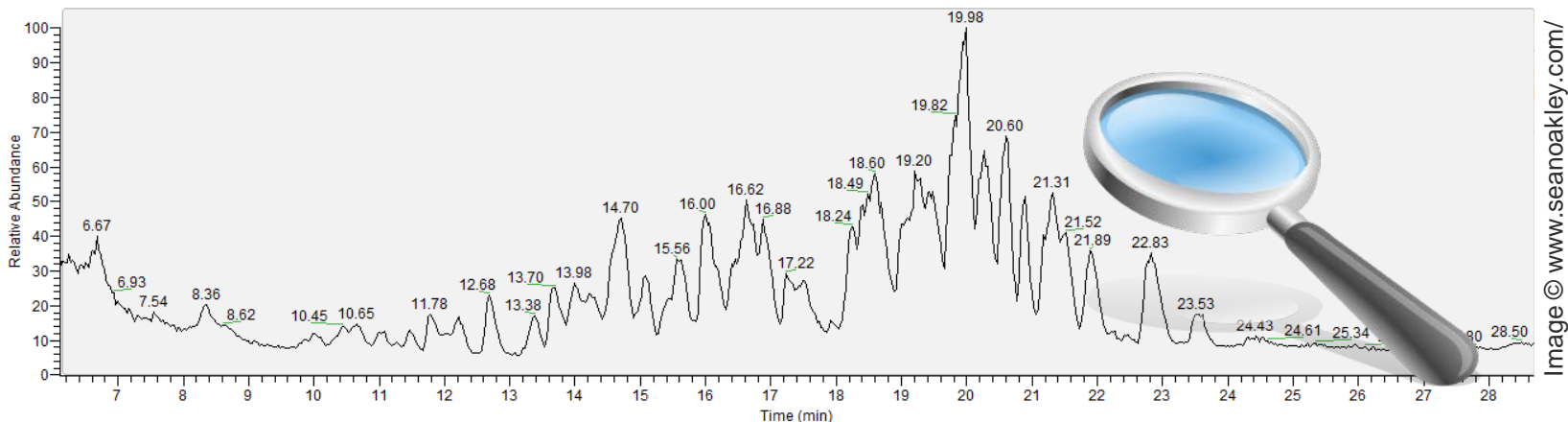


Image © www.seanoakley.com/

Assoc. Prof. Dr. Emma L. Schymanski

FNR ATTRACT Fellow and PI in Environmental Cheminformatics
Luxembourg Centre for Systems Biomedicine (LCSB), University of Luxembourg
Email: emma.schymanski@uni.lu

...and many colleagues who contributed to my science over the years!



Luxembourg



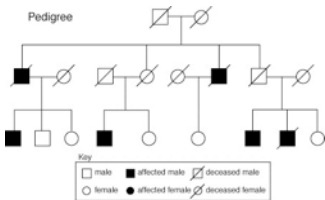
University of Luxembourg & LCSB



- Uni Lu was founded in 2003
 - Teenage years!
- LCSB was founded in 2009
 - Young and very dynamic working environment!



LCSB: Luxembourg Centre for Systems Biomedicine



Family Studies

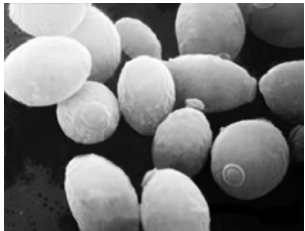
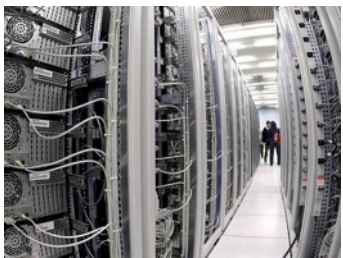
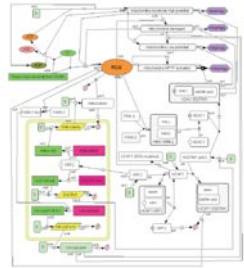


Patients



Longitudinal Cohorts

Pathway and Network Analysis
Computational Models
Machine Learning



Yeast



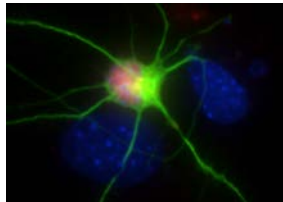
Zebrafish



Microbiome



Mouse









Human




LET'S MAKE IT HAPPEN


Members with access to **Environmental Cheminformatics**


-  **Adelene Lai** @adelene.lai
Given access 2 months ago
-  **Anjana Elapavalore** @anjana.elapavalore
Given access 4 weeks ago
-  **Corey Griffith** @corey.griffith
Given access 2 months ago
-  **Emma Schymanski** @emma.schymanski It's you
Given access 2 months ago
-  **German Andres Preciat Gonzales** @german.preciat
Given access 2 months ago


 **Hiba Hiba** @hiba.hiba
Given access 4 weeks ago

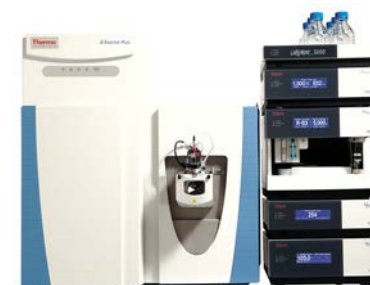
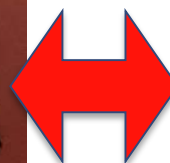
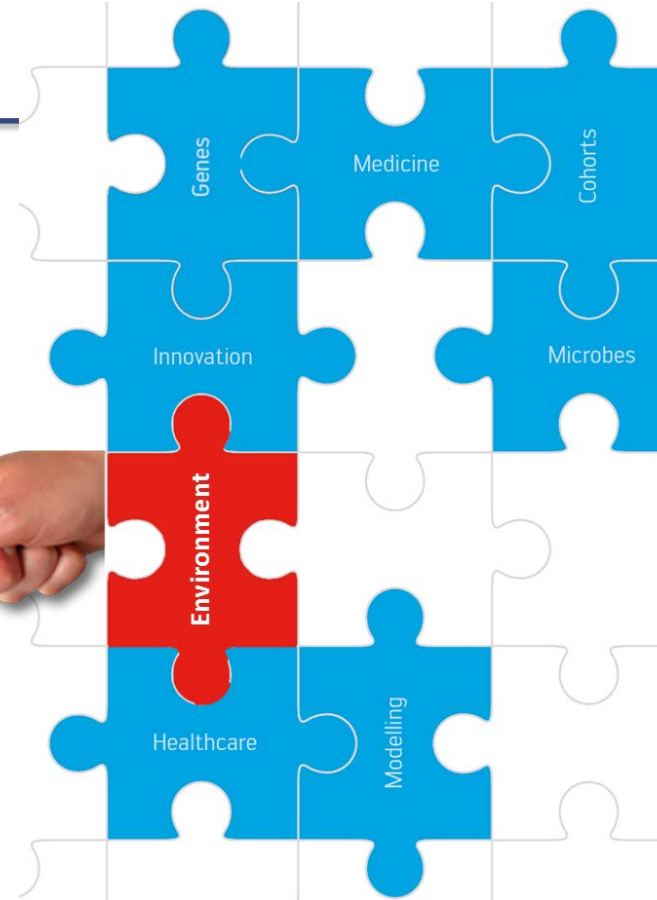
 **Jessy Krier** @jessy.krier
Given access 1 week ago

 **Lorenzo Favilli** @lorenzo.favilli
Given access 2 months ago

 **Mira Narayanan** @mira.narayanan
Given access 4 weeks ago

 **Randolph Singh** @randolph.singh
Given access 2 months ago

 **Todor Kondić** @todor.kondic
Given access 2 months ago

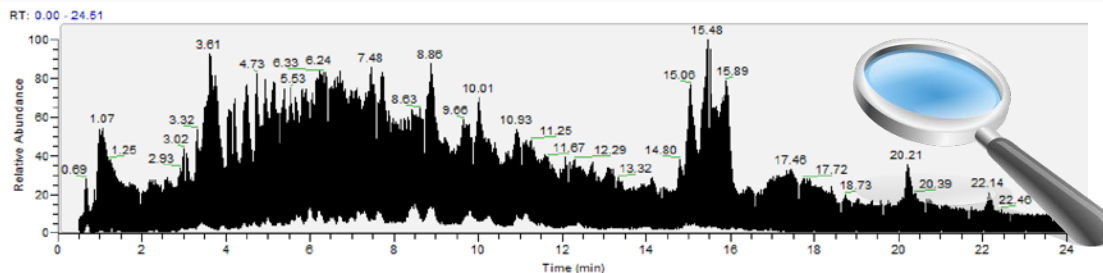


Luxembourg National
Research Fund

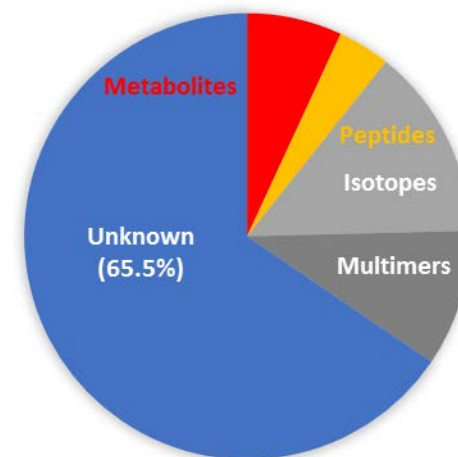
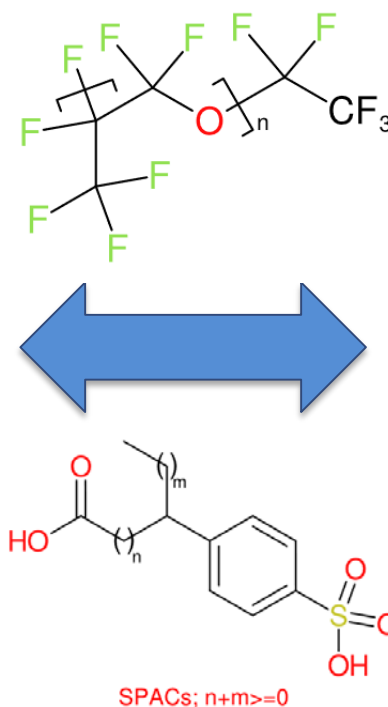
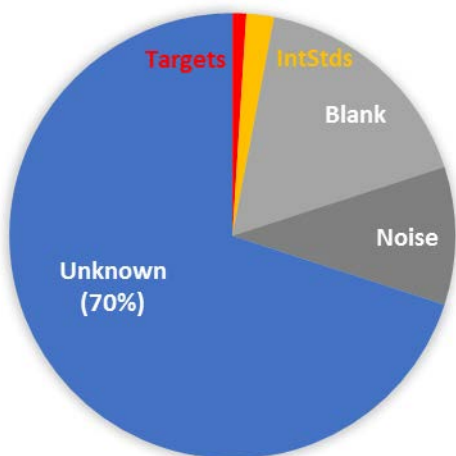


UNIVERSITÉ DU
LUXEMBOURG

Our challenge? We still have many unknowns ...

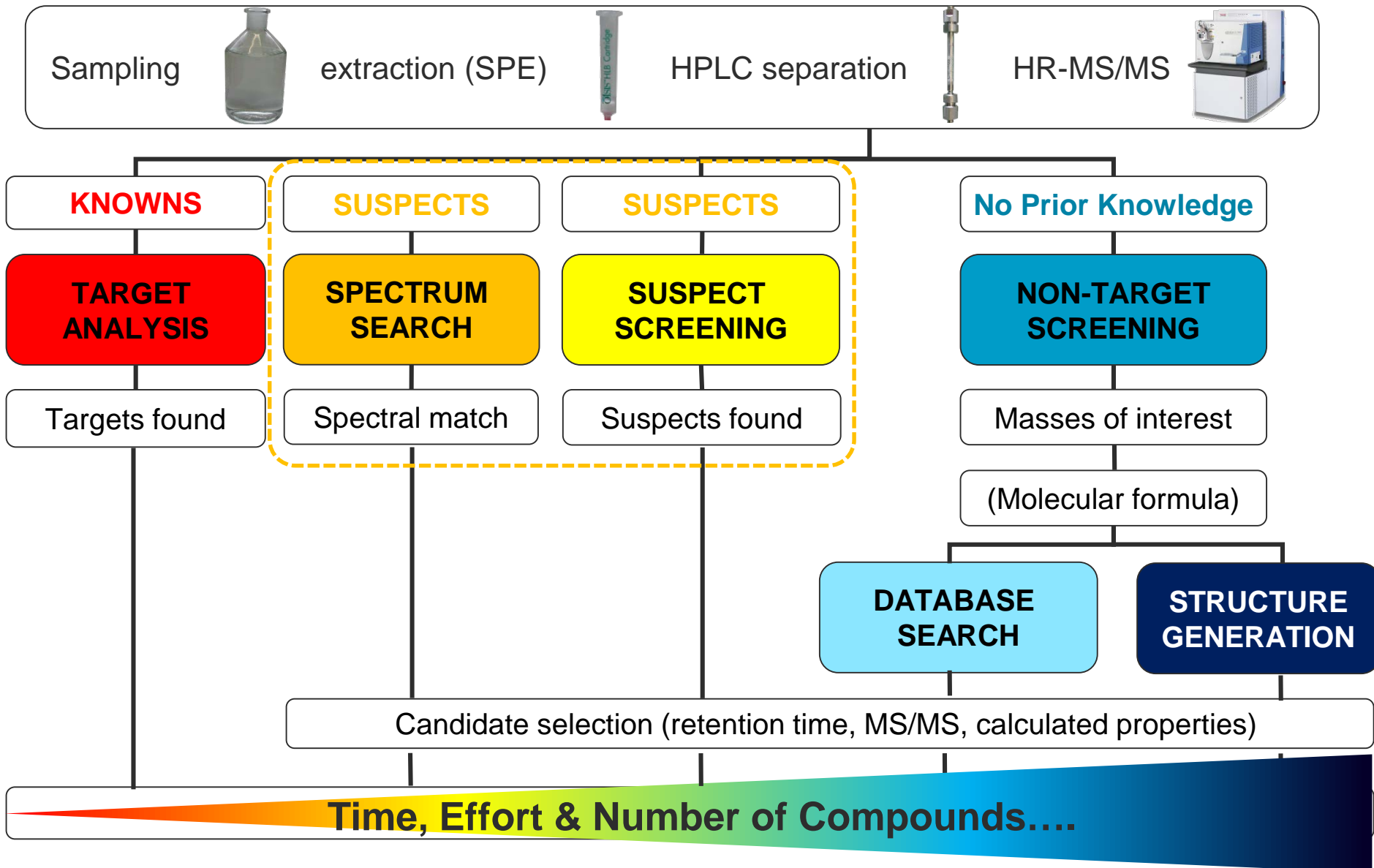


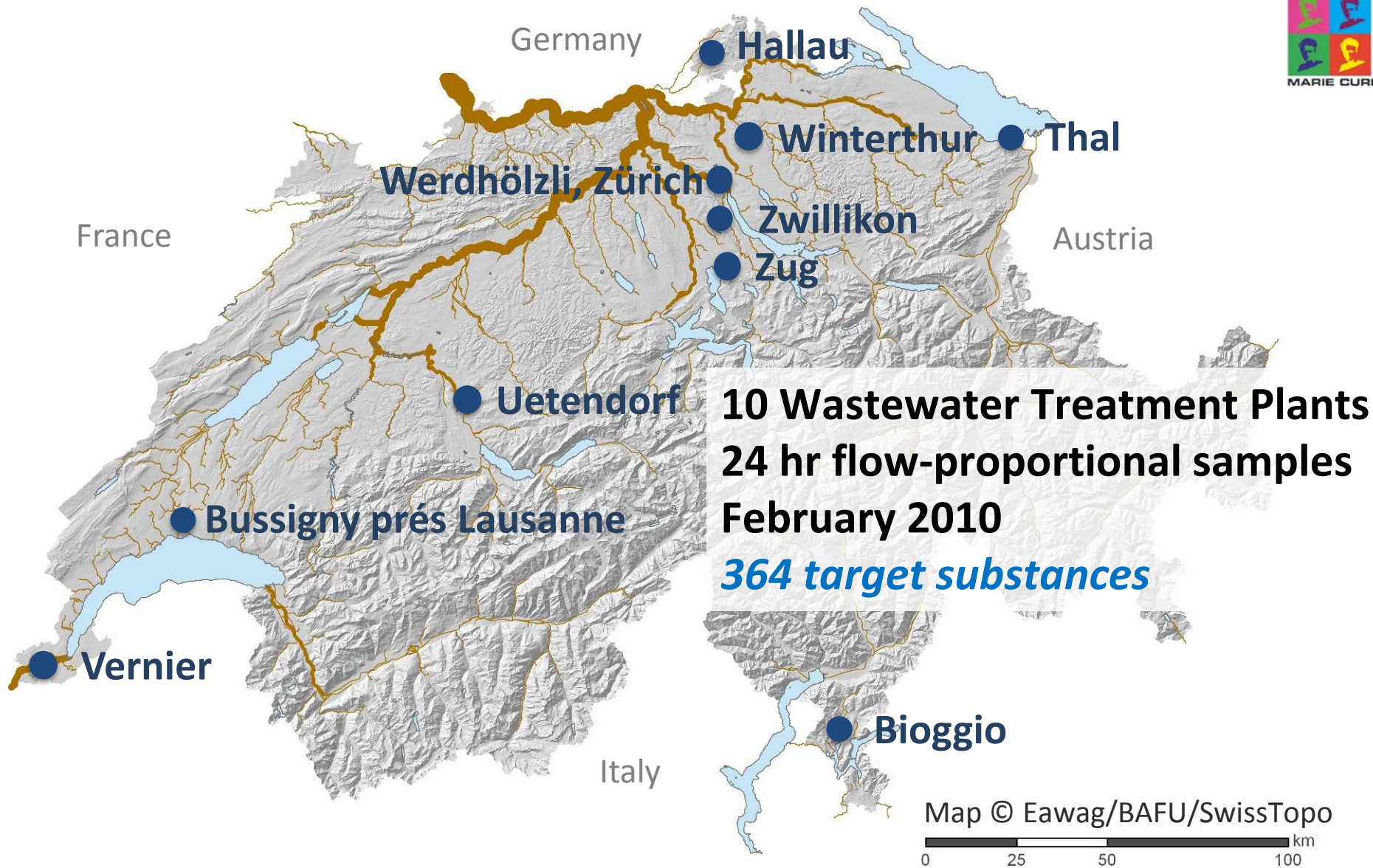
Wastewater



Cells

Target, Suspect and Non-Target Screening





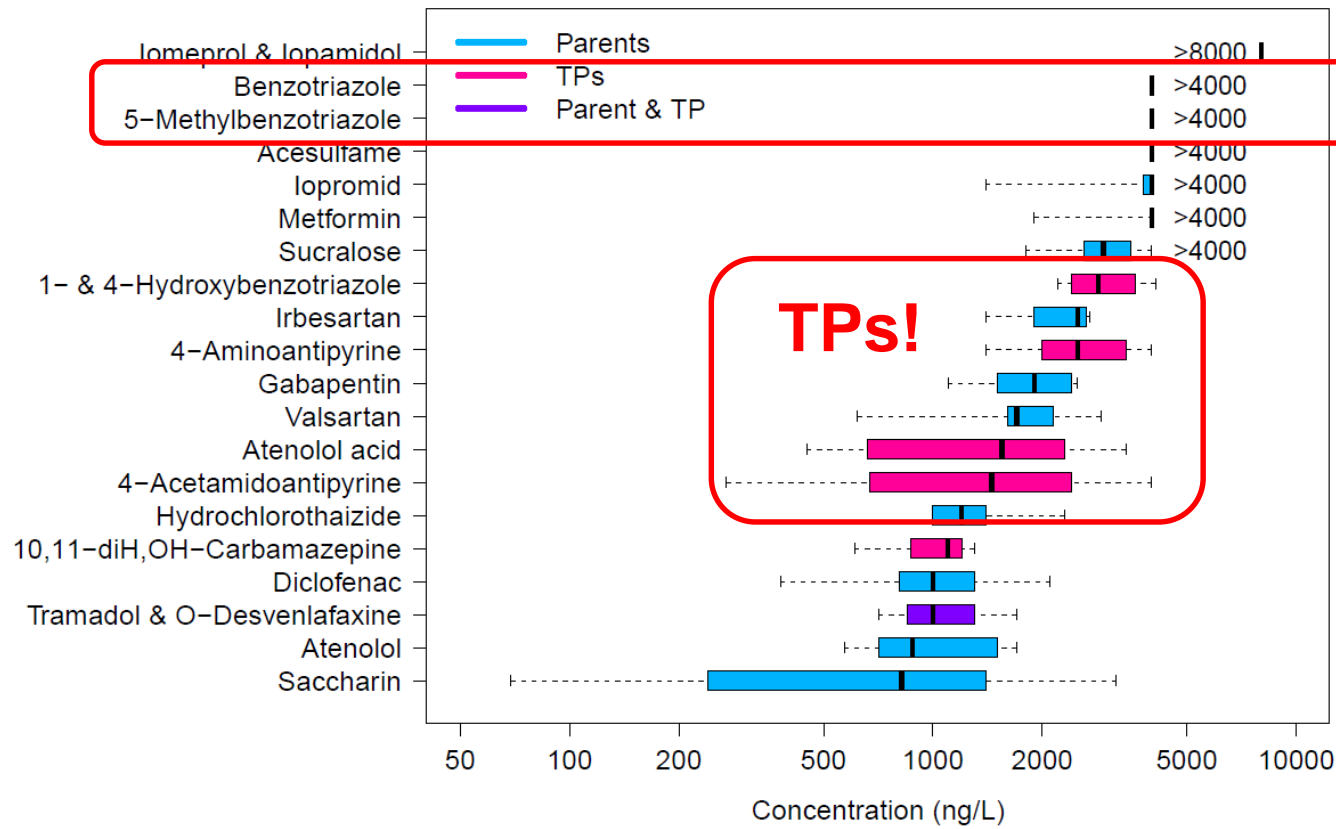
Target Analysis: Status Quo (>364 targets)



Target List

TARGET ANALYSIS

Targets found



Confirmation and quantification of compounds present

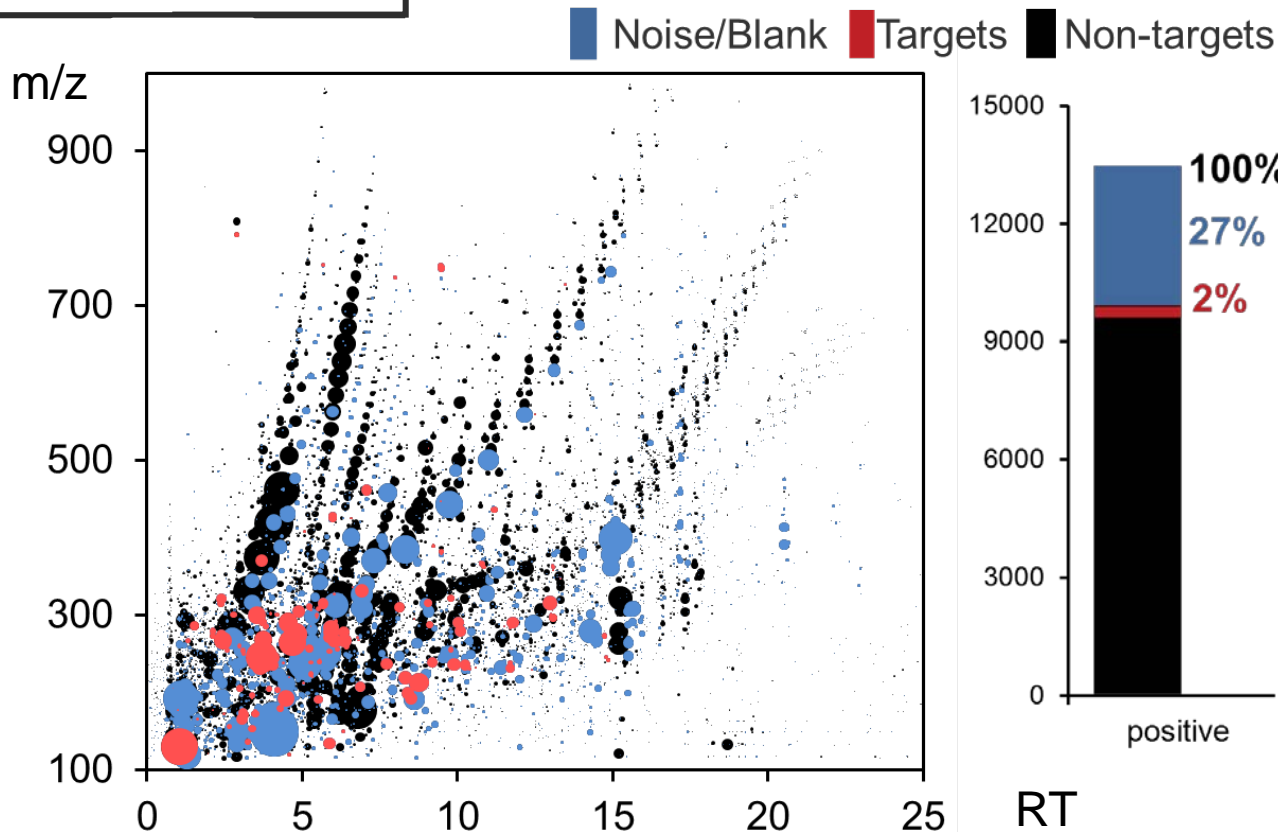
Target Analysis: Status Quo (>364 targets)



Target List

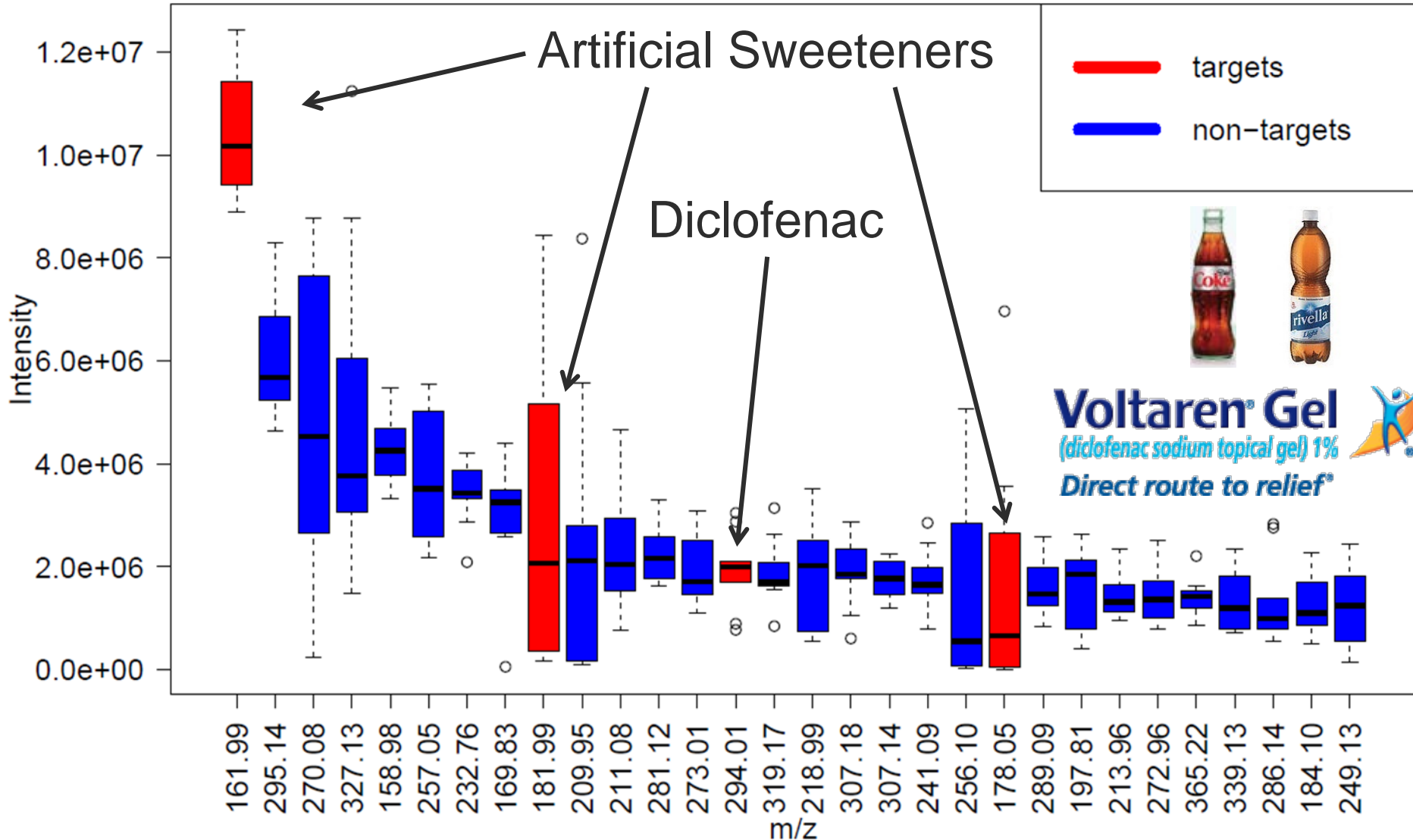
TARGET ANALYSIS

Targets found



Confirmation and quantification of compounds present

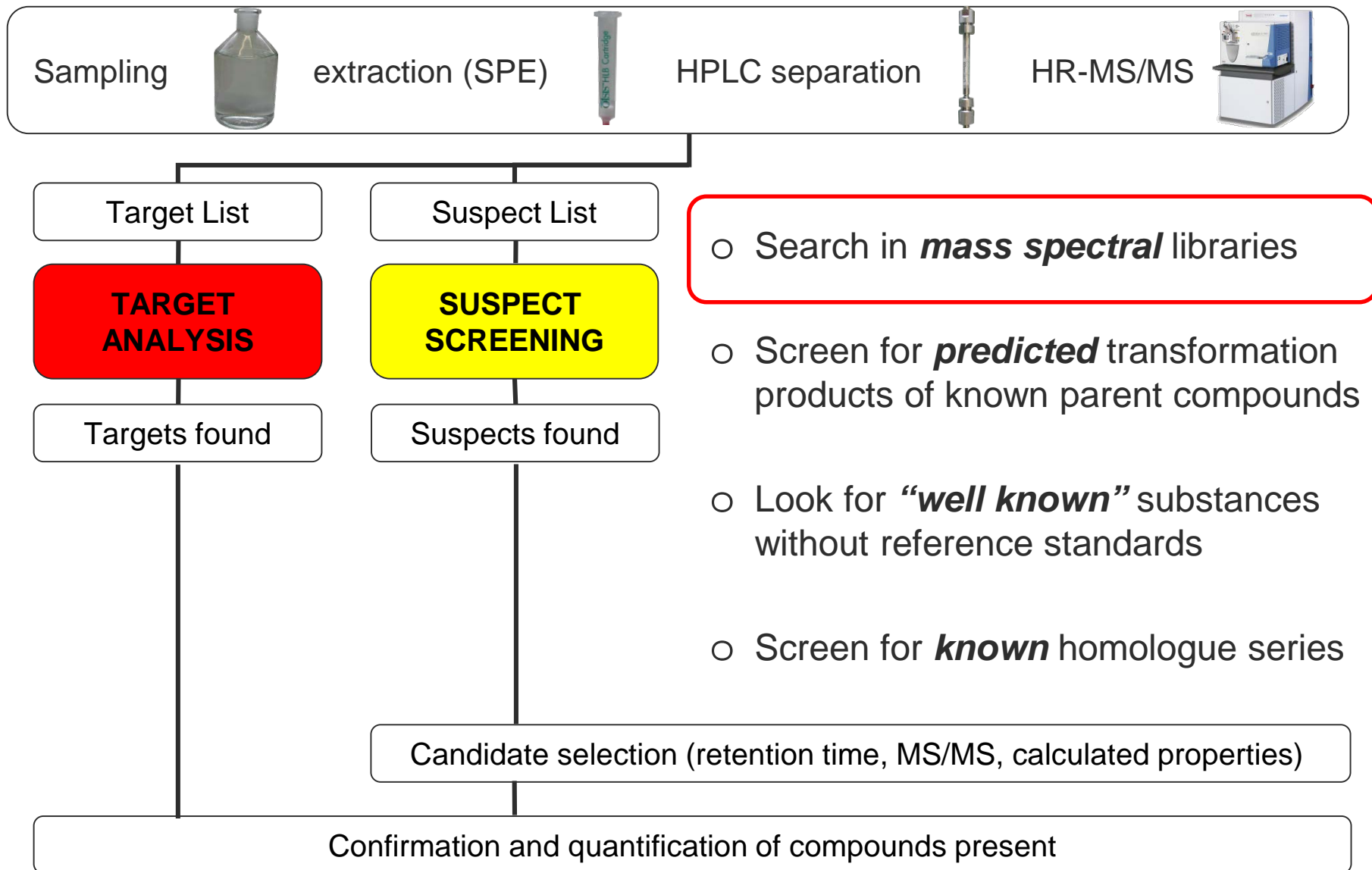
Swiss Wastewater: Top 30 Peaks (ESI-)



Pictures: www.coca-cola-com; www.rivella.ch; www.voltargengel.com

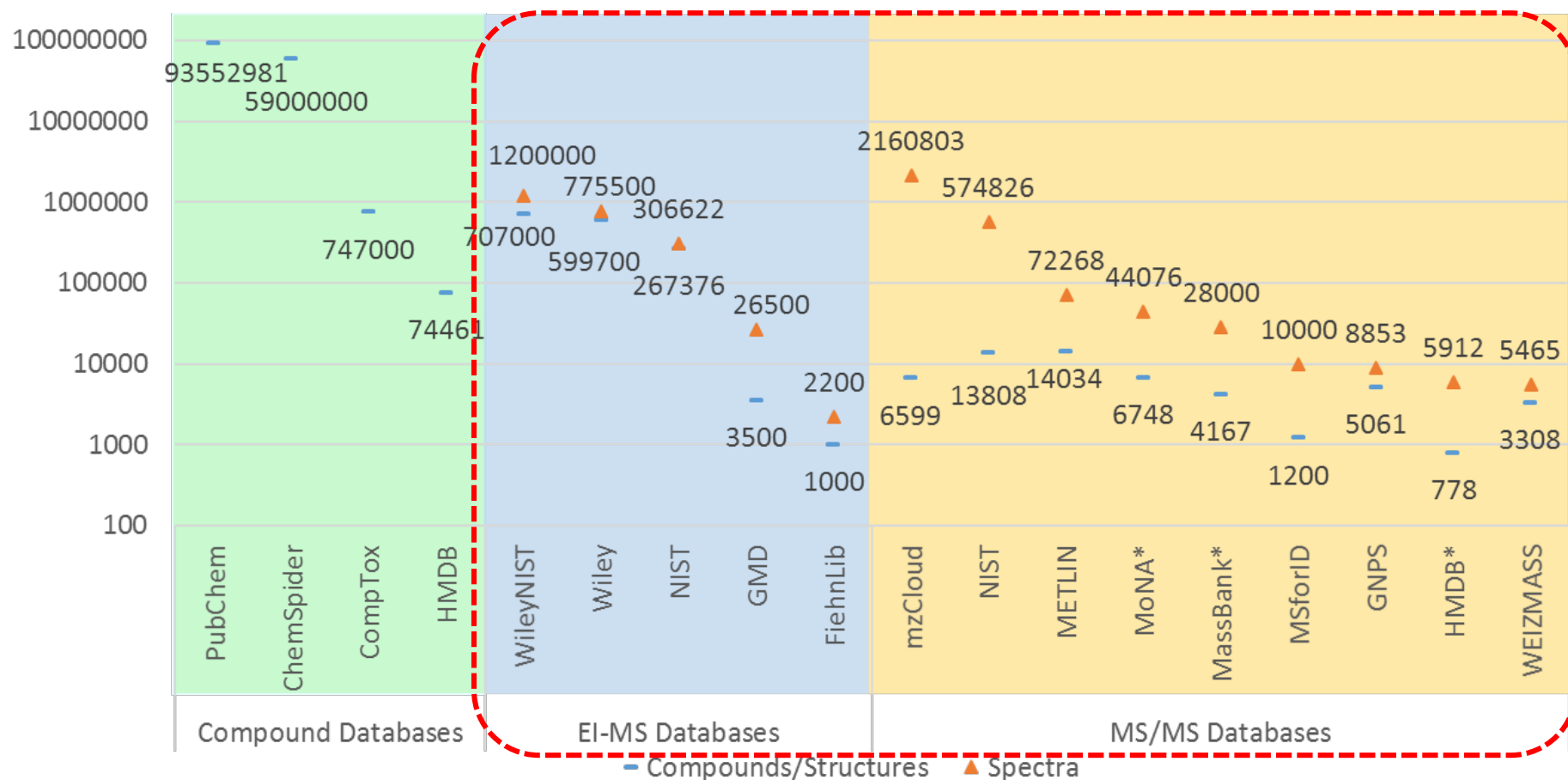
Schymanski et al. (2014), ES&T, 48: 1811-1818. DOI: [10.1021/es4044374](https://doi.org/10.1021/es4044374)

Suspect Screening: Different Approaches



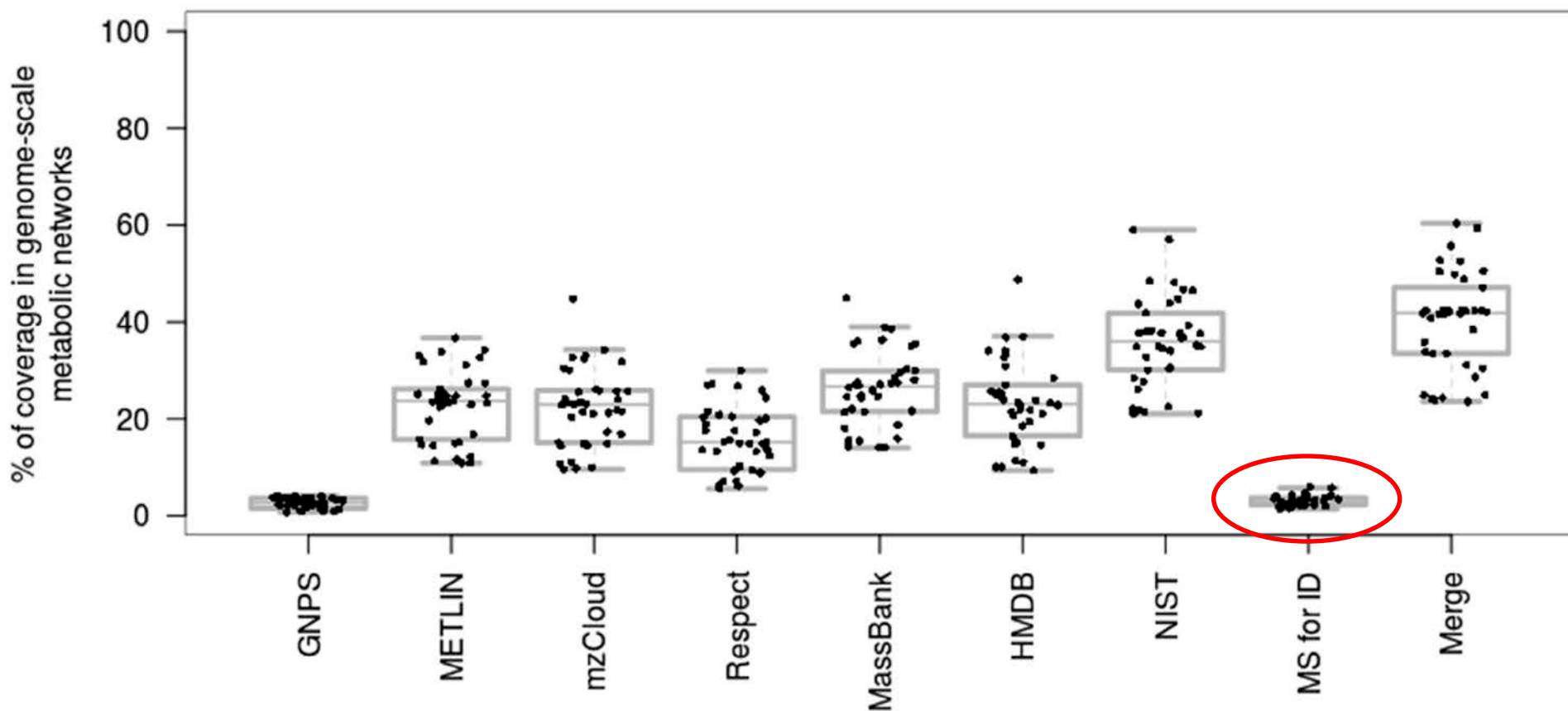
Searching Mass Spectral Libraries

o ... which one?



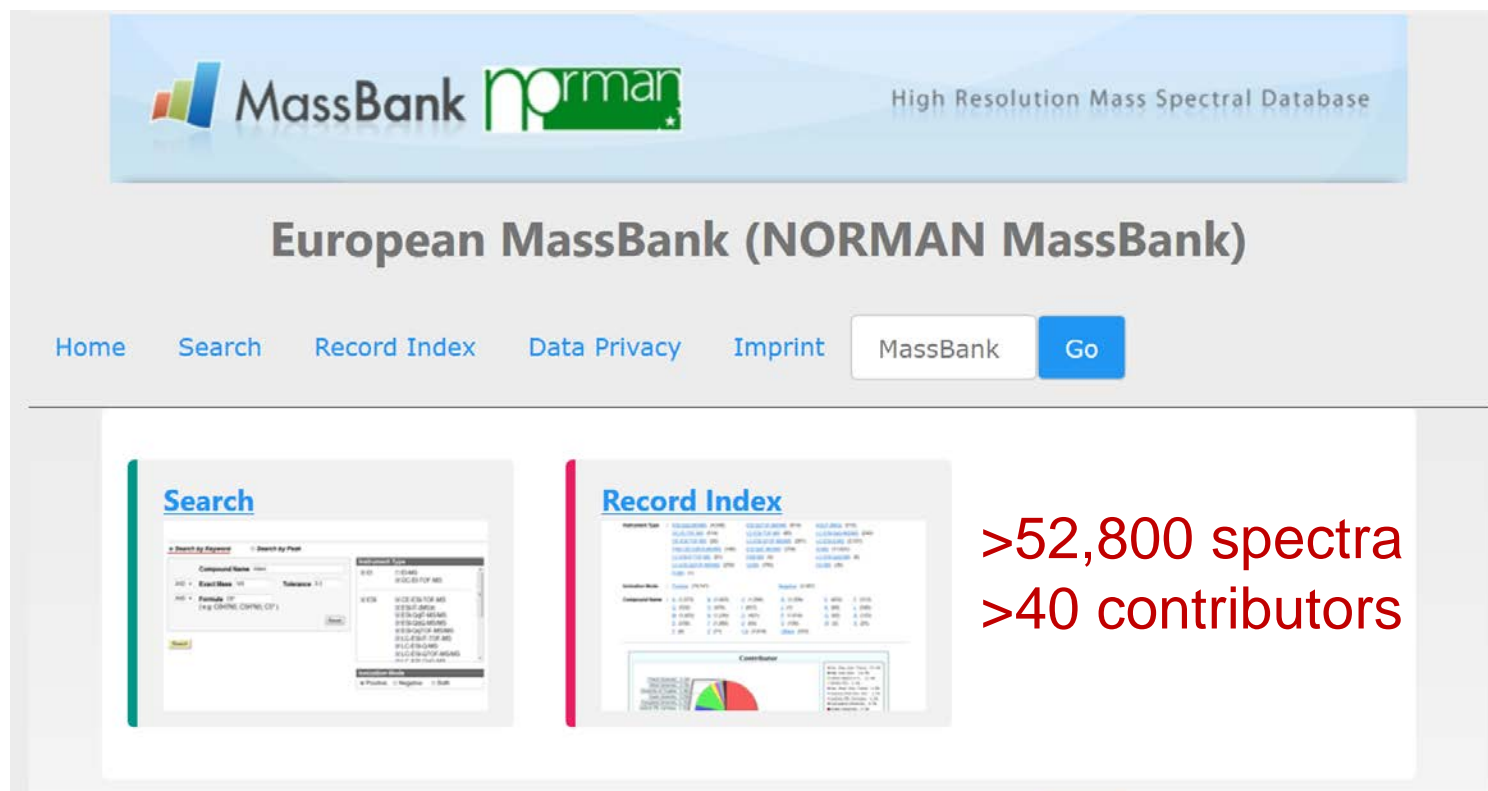
Mind the Gap!

- Best library to choose depends highly on your dataset
 - Example: MSforID (<https://msforid.com/>) is poor for metabolic networks – but great for forensic toxicology!



<http://massbank.eu/MassBank>

<https://github.com/MassBank/MassBank-data/>

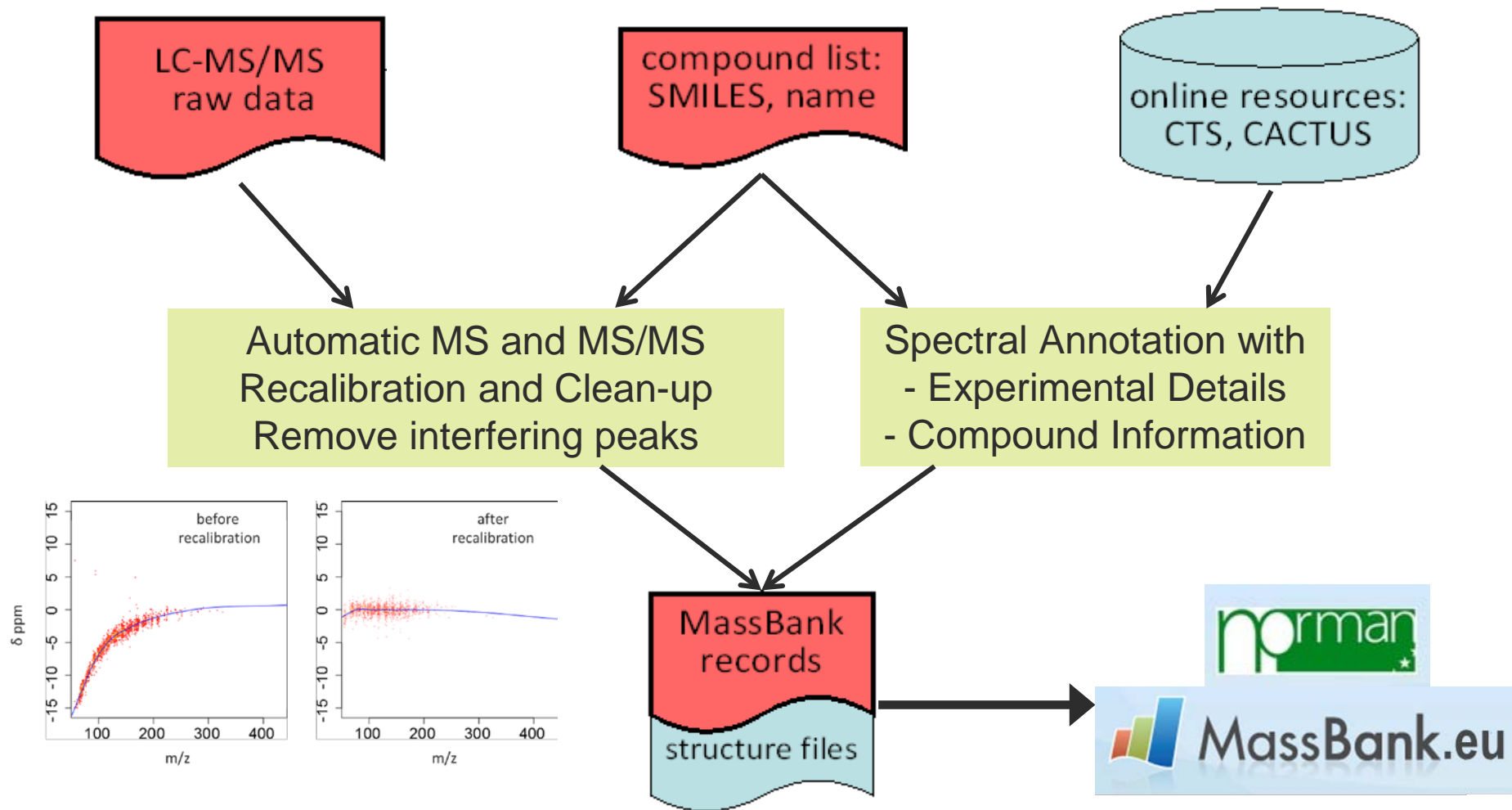


The screenshot displays the MassBank website header with the logo and the text "High Resolution Mass Spectral Database". Below the header is the title "European MassBank (NORMAN MassBank)". A navigation menu includes "Home", "Search", "Record Index", "Data Privacy", and "Imprint". A search bar contains the text "MassBank" and a "Go" button. Below the navigation are two preview images: "Search" showing a search interface with a list of results, and "Record Index" showing a table of records and a pie chart labeled "Contributor".

>52,800 spectra
>40 contributors

MassBank-data validation status

build **passing**

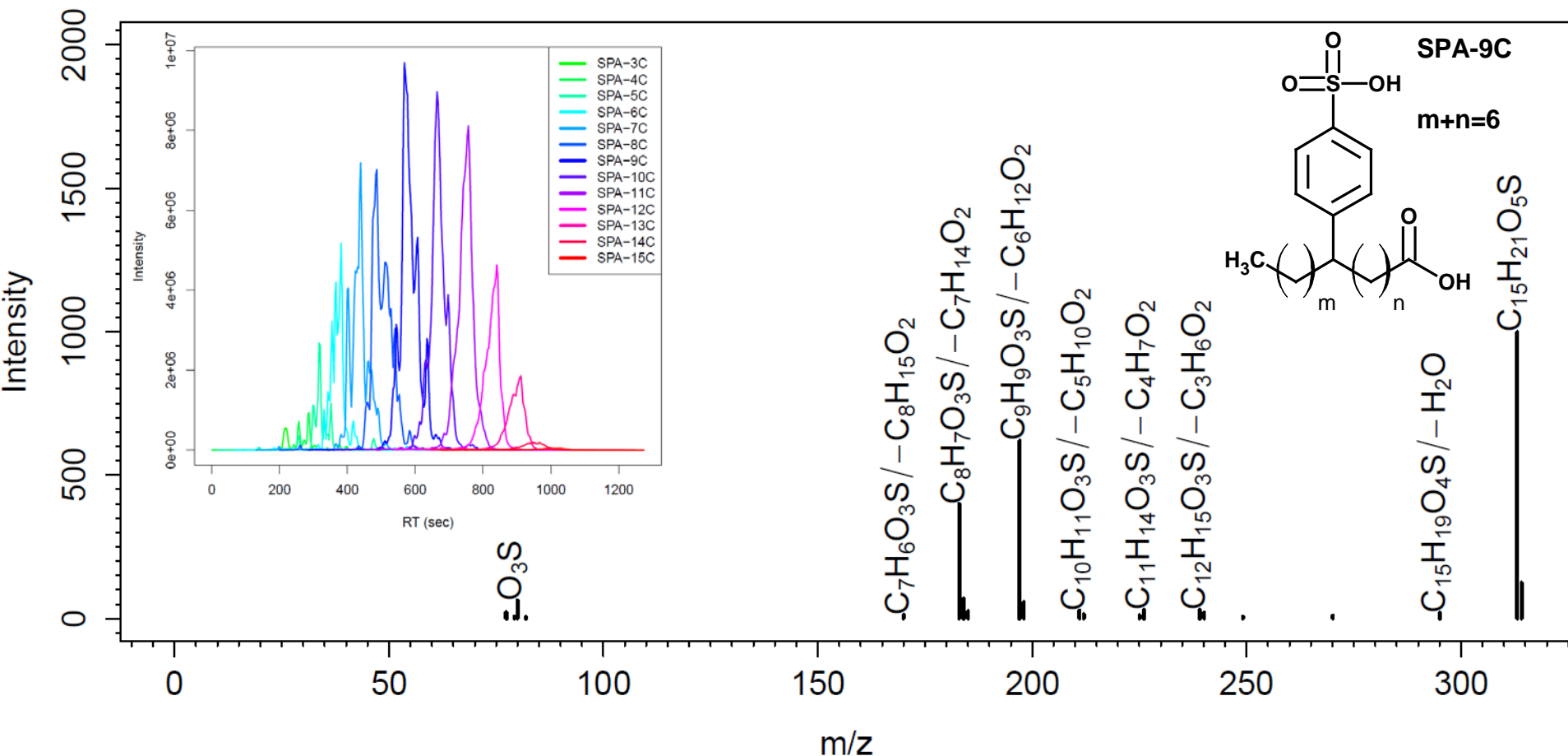


<https://github.com/MassBank/RMassBank/>
<http://bioconductor.org/packages/RMassBank/>

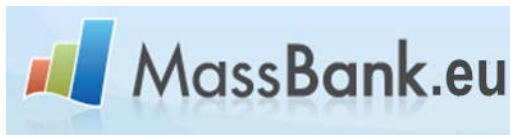
Chromatography and MS/MS Annotation

<https://github.com/MassBank/RMassBank/>

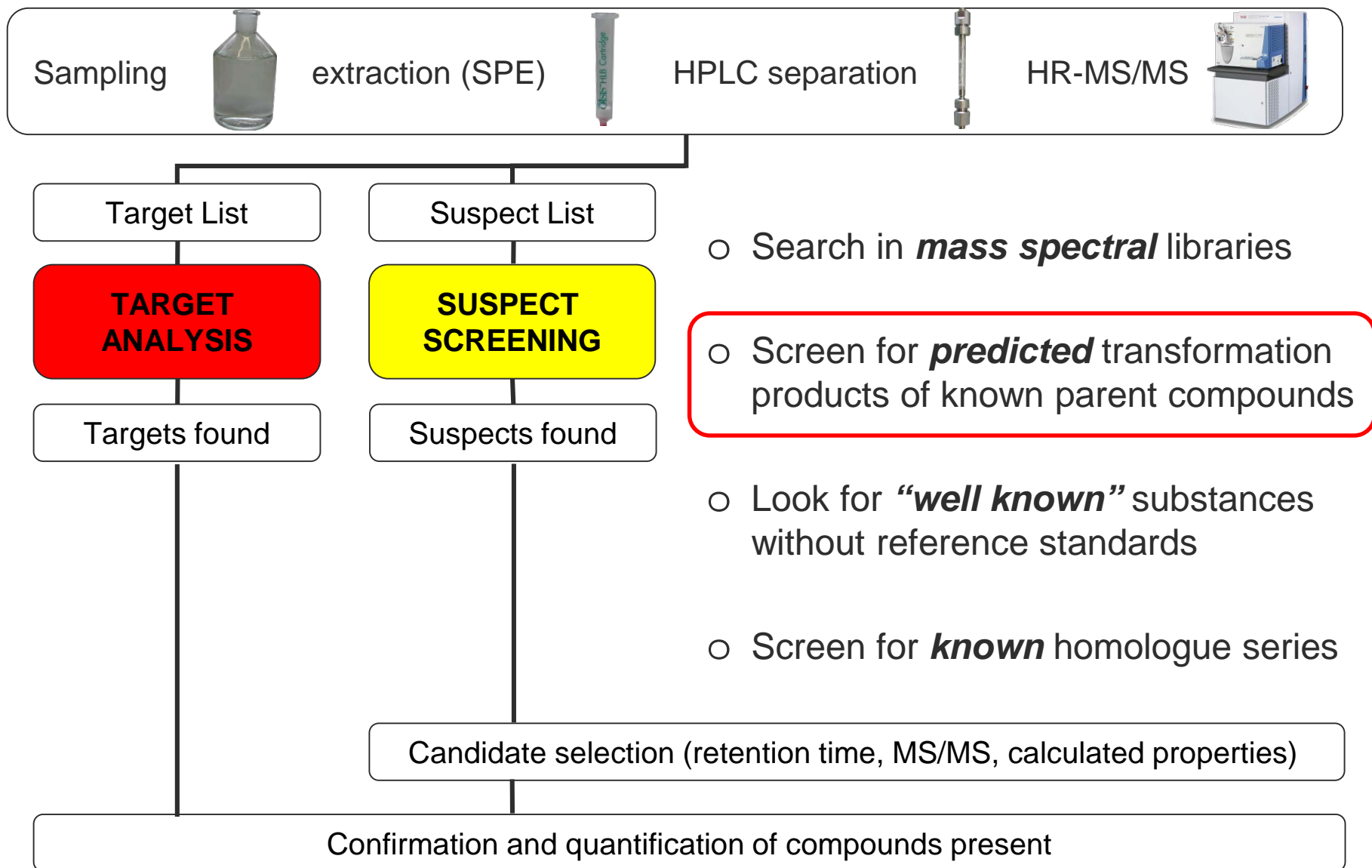
RMassBank



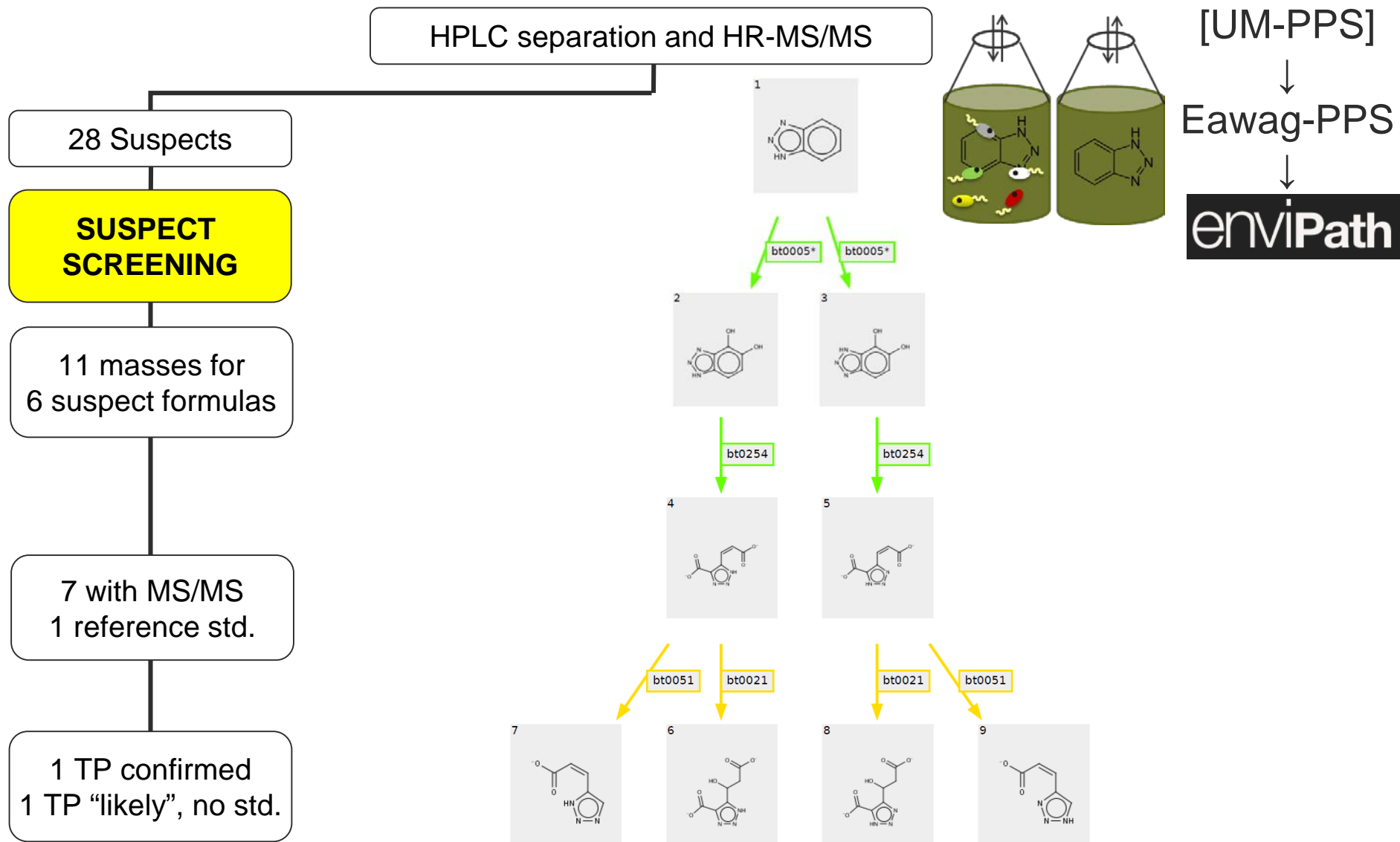
Formulas: <http://sourceforge.net/projects/genform/>
Meringer *et al.*, 2011, *MATCH* 65, 259-290
Data: Schymanski *et al.*, 2014, *ES&T*, 48:
1811-1818. DOI: 10.1021/es4044374



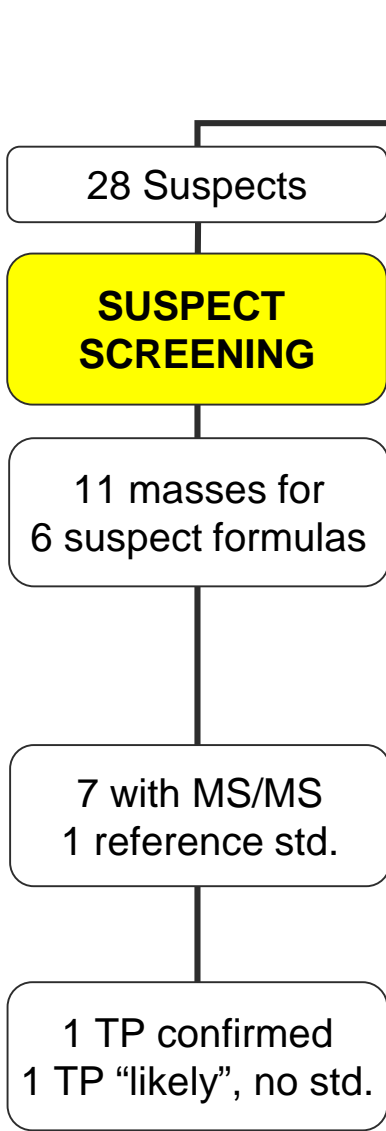
Literature: LIT00034,35
Sample: ETS00002
Standard: ETS00016,17,19,20



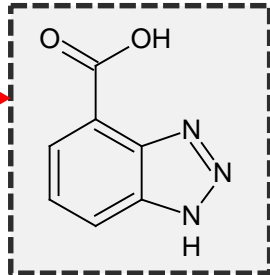
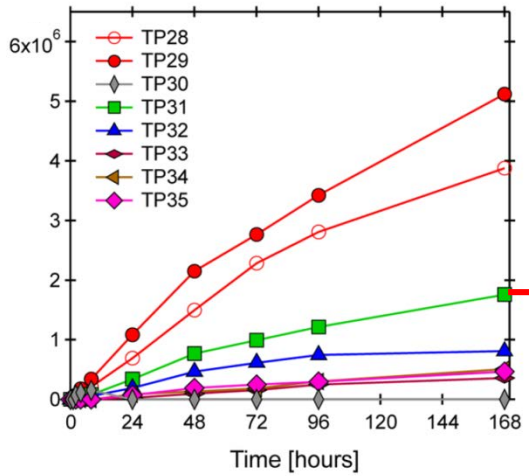
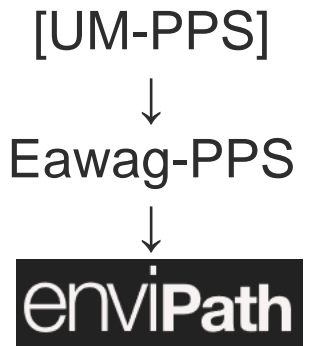
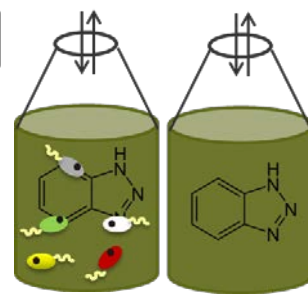
Suspect Screening: Benzotriazole TPs



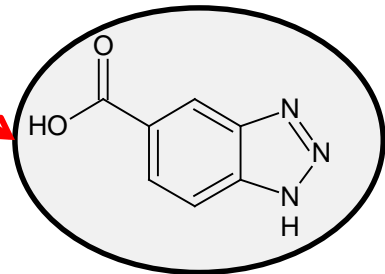
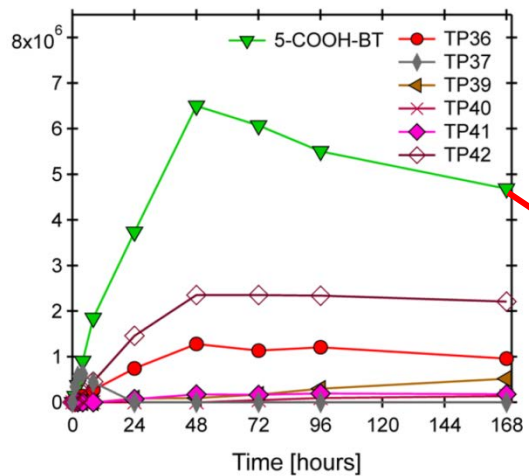
Suspect Screening: Benzotriazole TPs



HPLC separation and HR-MS/MS

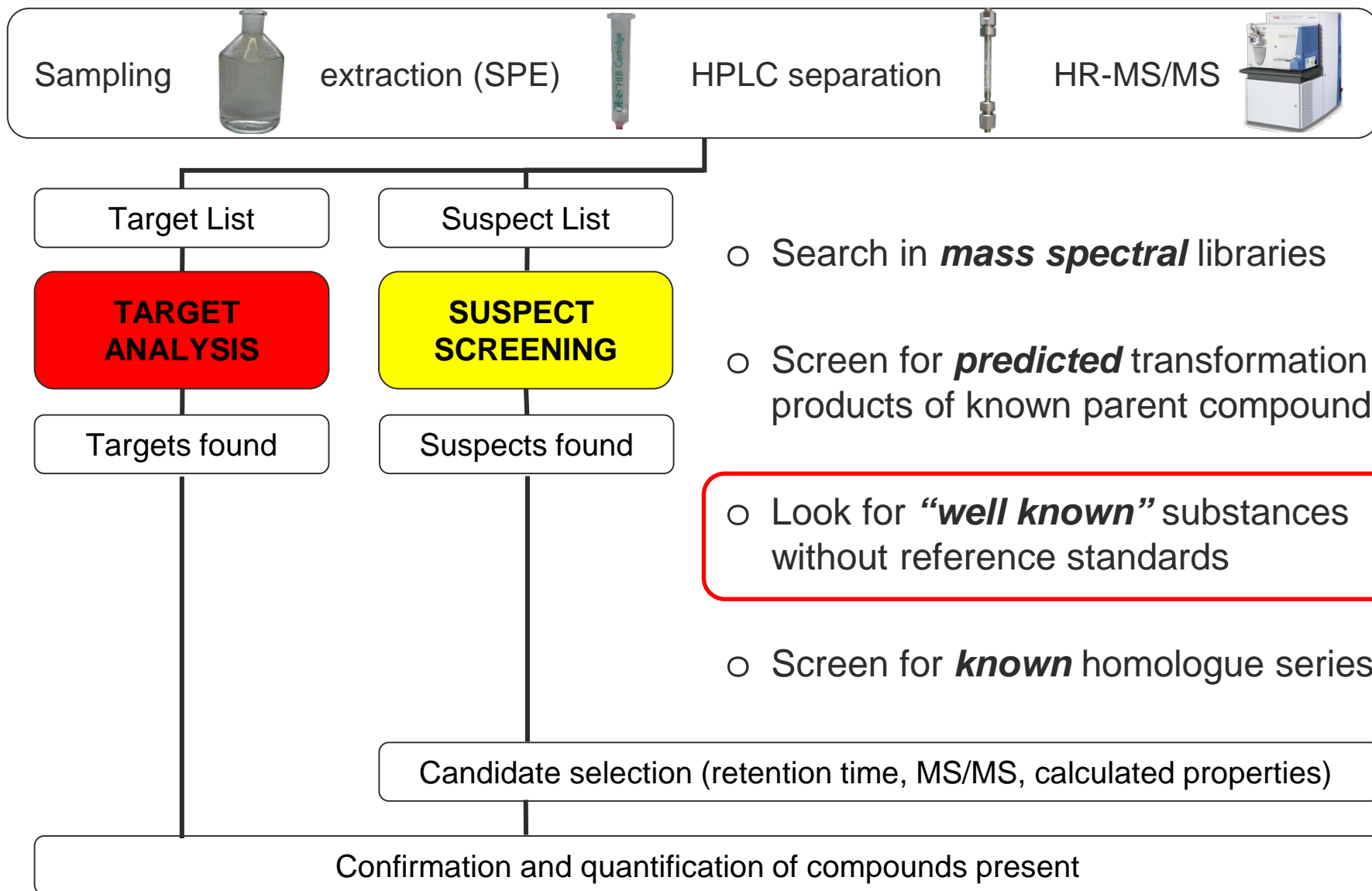


- Predicted with Eawag-PPS
- No standard
- Not in databases (at that time)



- Confirmed with reference std.
- Observed in WWTP effluents

Suspect Screening: Different Approaches



(European) Environmental Community (subset!)





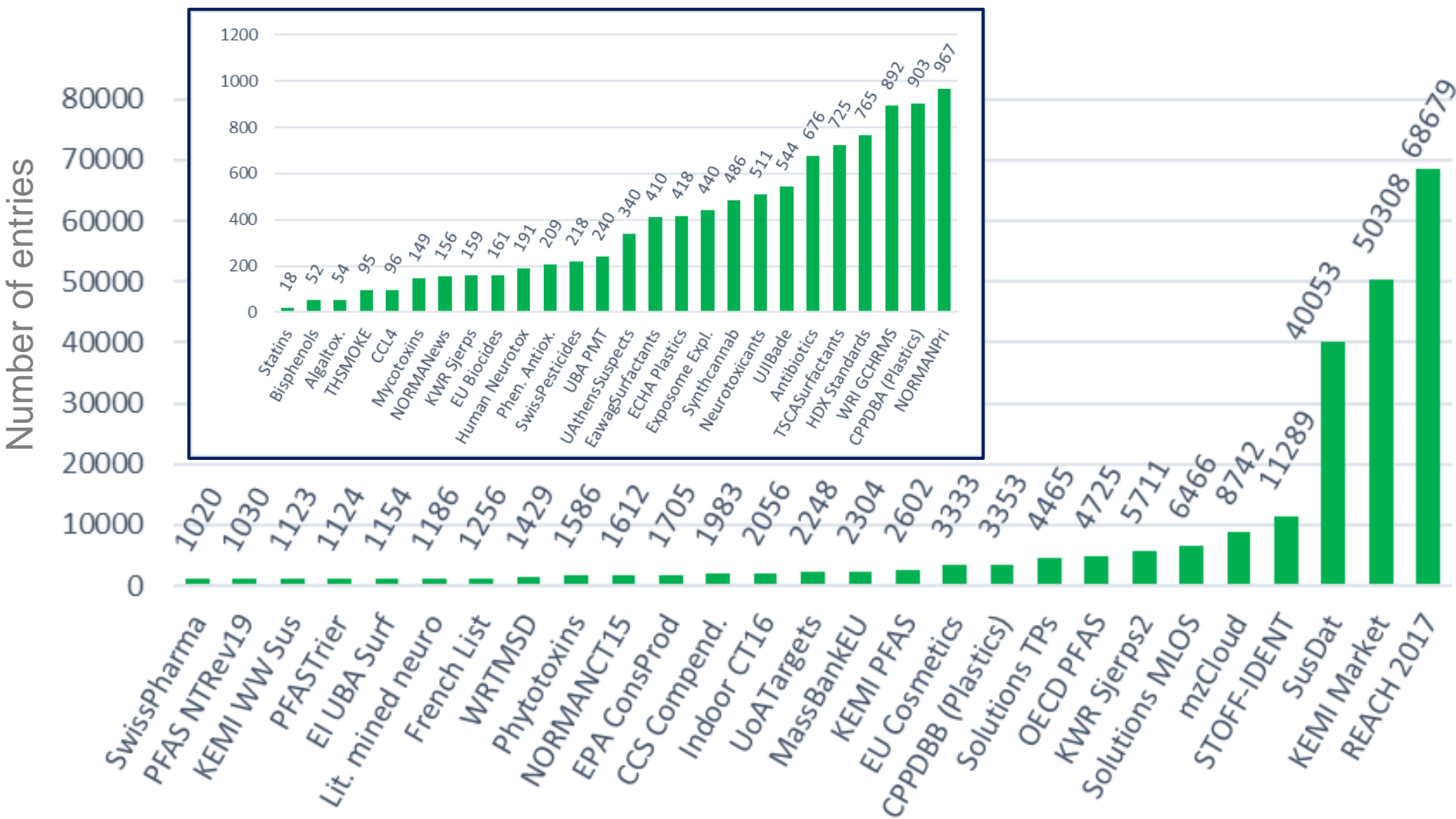
NORMAN Suspect List Exchange – NORMAN SLE

<https://www.norman-network.com/nds/SLE/>

No.	Abbreviation	Description	Link to full list	Link to InChIKey list	References
S0	SUSDAT	Merged NORMAN Suspect List: SusDat	Interactive Data table (updating...) CompTox SUSDAT List	MS-ready InChIKeys (1/03/2018)	A merged list of >40,000 structures from suspect lists. See interactive version . Compiled by Reza Aalizadeh, University of Athens, including RTI and toxicity values, support by Nikiforos Alygizakis, EI. <i>Work in progress ... please report any issues!</i> DOI: 10.5281/zenodo.2664077
S1	MASSBANK	NORMAN Compounds in MassBank	CSV, XLSX with Fragments (3/10/2017) CompTox MassBank EU Reference List CompTox MassBank EU Special Cases CompTox Fragment Download	MassBankEUInChIKeys (17/06/2019)	www.massbank.eu Stravs <i>et al.</i> 2013. DOI: 10.1002/jms.3131 DOI: 10.5281/zenodo.2621390
S2	STOFFIDENT	HSWT/LFU STOFF-IDENT Database of Water-Relevant Substances	STOFF-IDENT Contents (6/09/2017) CompTox STOFF-IDENT List	STOFF-IDENT InChIKeys (6/09/2017)	The database enables the search for exact masses from target or unknown lists and the automatic use of a Retention Time Index. See: https://www.lfu.bayern.de/stoffident/#!home (single search for free; batch search after free registration). DOI: 10.5281/zenodo.2621451
S3	NORMANCT15	NORMAN Collaborative Trial Targets and Suspects	LC-MS: CSV, XLSX (3/10/2017) GC-MS: CSV, XLSX (3/10/2017) CompTox NORMANCT15 List	LC-MS InChIKeys (31/10/2016) GC-MS InChIKeys	Schymanski <i>et al.</i> 2015. DOI: 10.1007/s00216-015-8681-7 DOI: 10.5281/zenodo.2621478

Suspect Screening

- [NORMAN Suspect List Exchange \(NORMAN-SLE\)](#): 58 lists!



<https://www.norman-network.com/nds/susdat/>

https://comptox.epa.gov/dashboard/chemical_lists/susdat



NORMAN SusDat: Suspect List Exchange Data Table

Structure & Properties

If no criteria is selected, the result of search will be the overall database.

Search criteria - Individual substance(s)

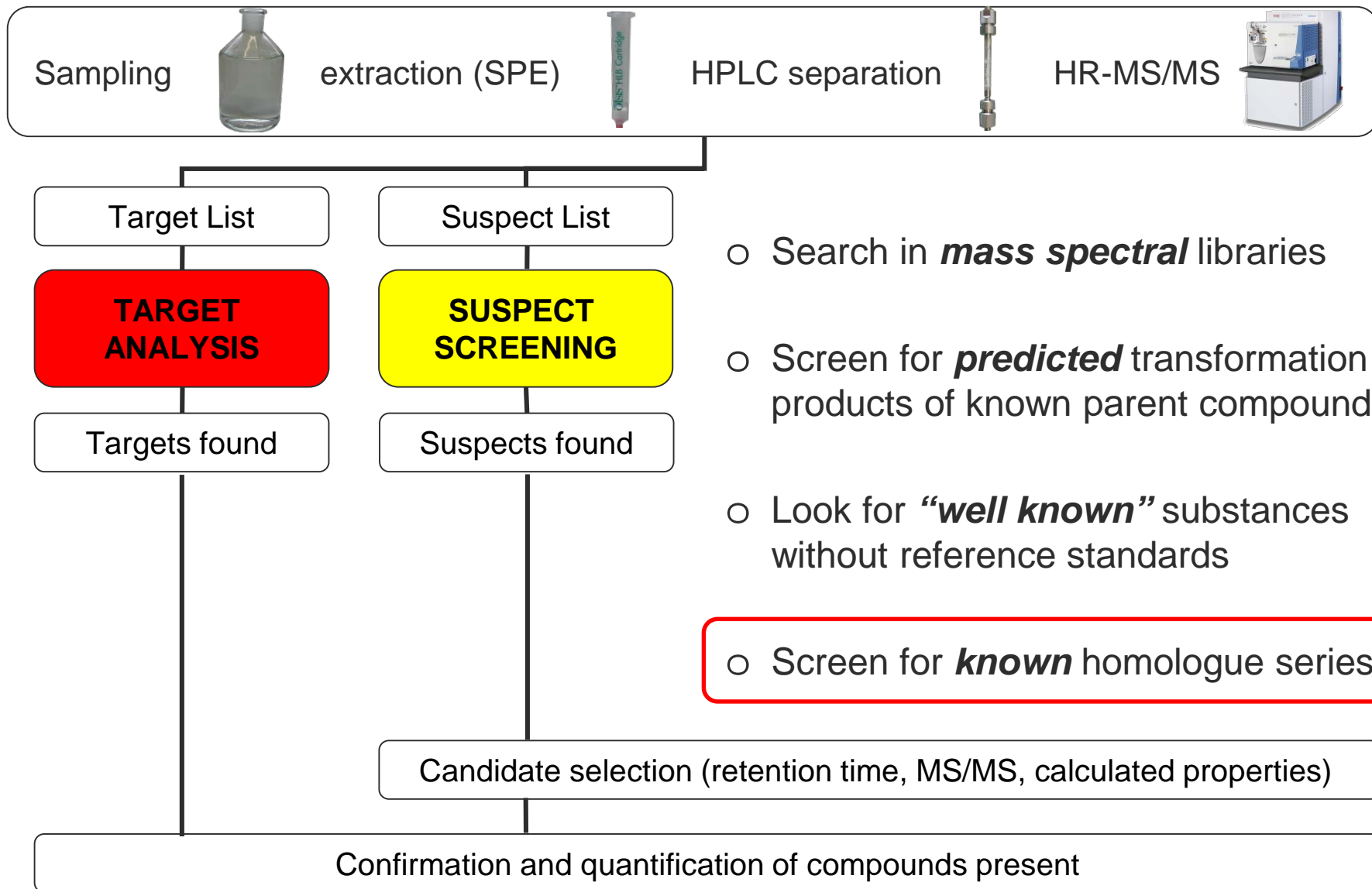
Search criteria - Batch search

S1/MASSBANK NORMAN Compounds in MassBank

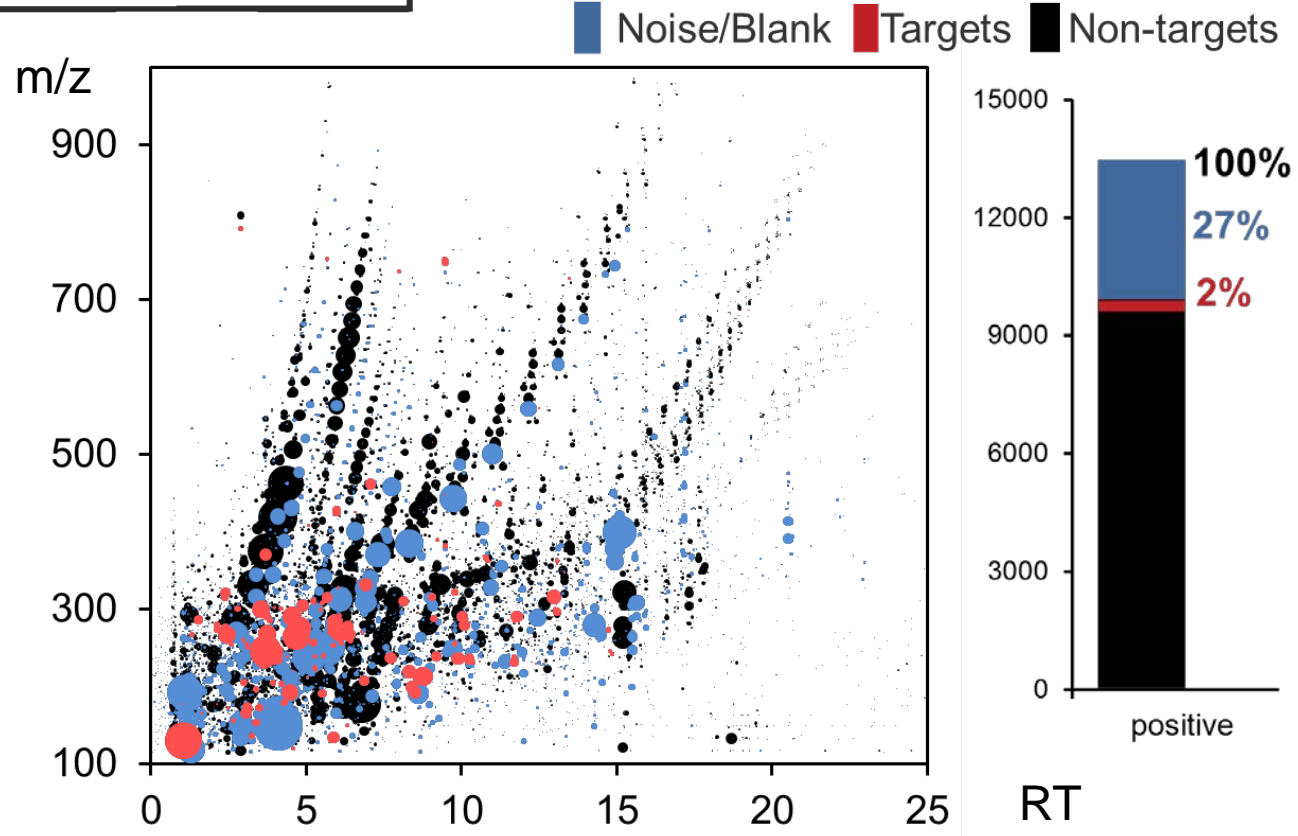
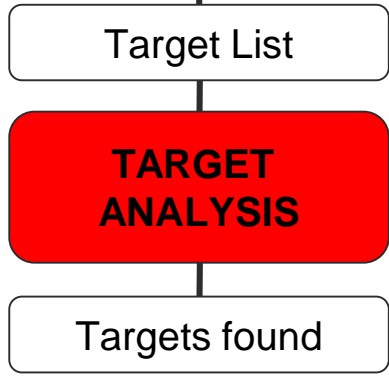
S2/STOFFIDENT HSWT/LfU STOFF-IDENT Database of Water-Relevant Substances

	Norman SusDat ID	Name	CAS_RN	Validation Level	SMILES
<input type="button" value="Reset"/>	<input type="text"/>	<input type="text"/>	<input type="text"/>	<input type="text"/>	<input type="text"/>
<input type="button" value="🔍"/>	NS00000001	Sulfaclozine	102-65-8	Level 4	<chem>c1cc(ccc1N)S(=O)(=O)Nc2cncc(n2)Cl</chem>
<input type="button" value="🔍"/>	NS00000002	Sulfachlorpyridazine	80-32-0	Level 2	<chem>c1cc(ccc1N)S(=O)(=O)Nc2ccc(nn2)Cl</chem>
<input type="button" value="🔍"/>	NS00000003	Sulfaguanidine	57-67-0	Level 2	<chem>c1cc(ccc1N)S(=O)(=O)NC(=N)N</chem>
<input type="button" value="🔍"/>	NS00000004	Sulfamerazine	127-79-7	Level 2	<chem>Cc1ccnc(n1)NS(=O)(=O)c2ccc(cc2)N</chem>
<input type="button" value="🔍"/>	NS00000005	Sulfamethizole	144-82-1	Level 2	<chem>Cc1nnc(s1)NS(=O)(=O)c2ccc(cc2)N</chem>
<input type="button" value="🔍"/>	NS00000006	Sulfamoxole	729-99-7	Level 2	<chem>Cc1c(cc(n1)NS(=O)(=O)c2ccc(cc2)N)C</chem>
<input type="button" value="🔍"/>	NS00000007	Sulfanilamide	1337-39-9	Level 4	<chem>c1cc(ccc1N)S(=O)(=O)N</chem>

Suspect Screening: Different Approaches



RECAP: Target Analysis: Status Quo (>364 targets)



Confirmation and quantification of compounds present

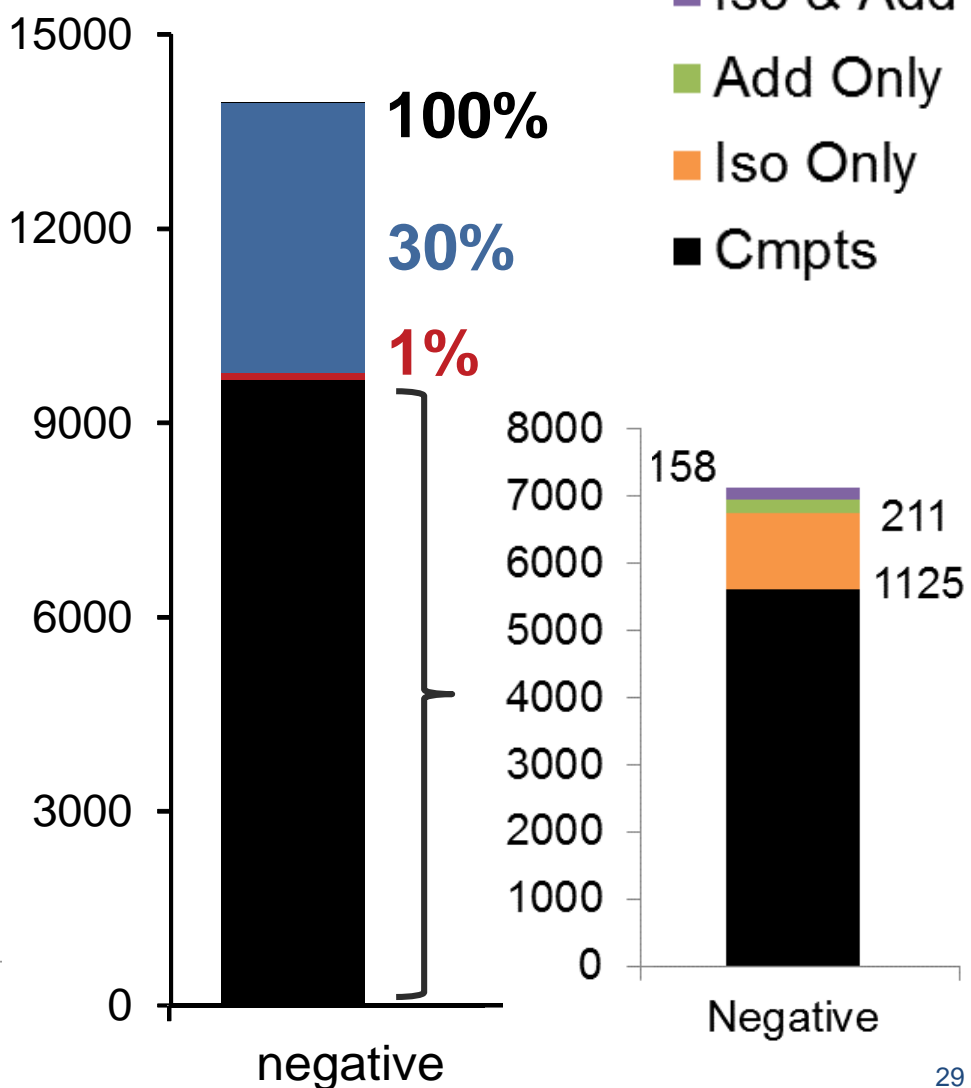
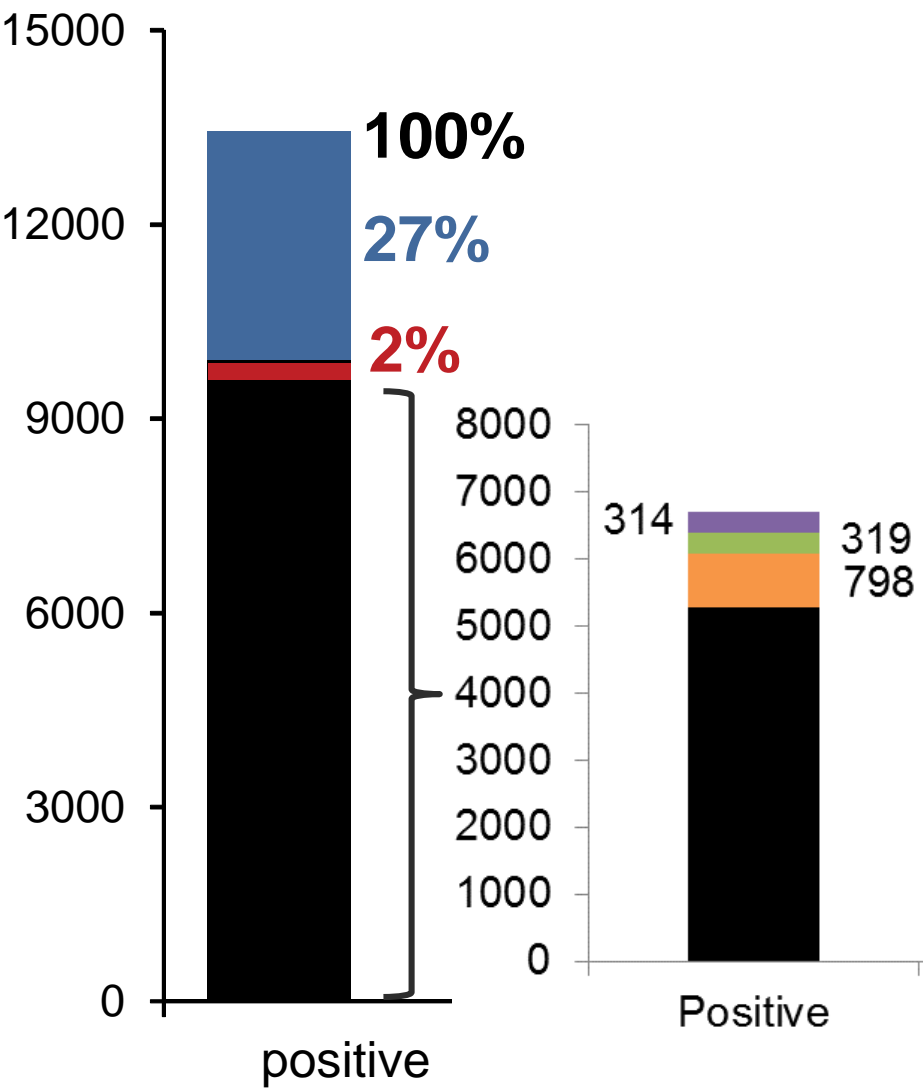
Grouping Isotopes and Adducts

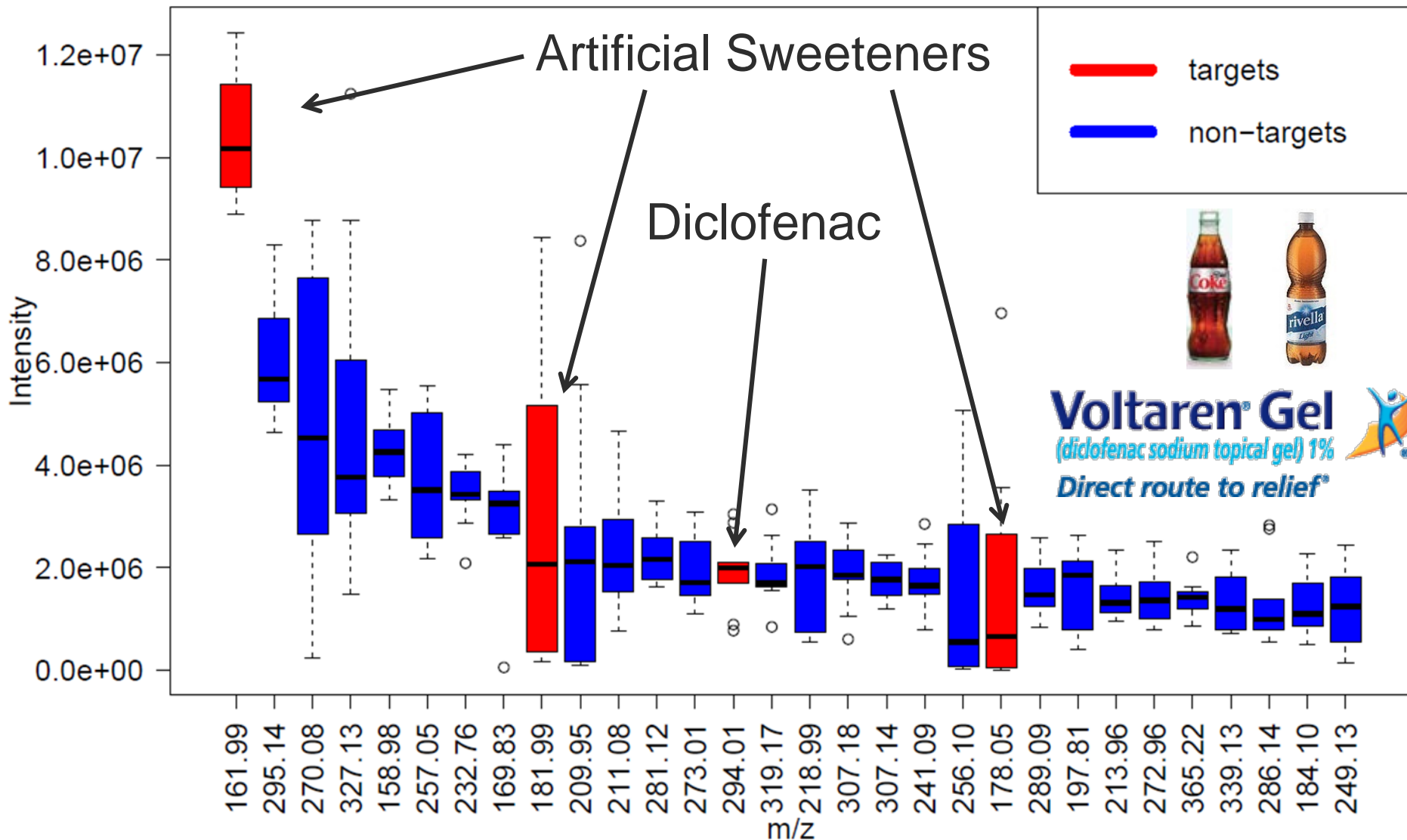


■ Noise/Blank ■ Targets ■ Non-targets

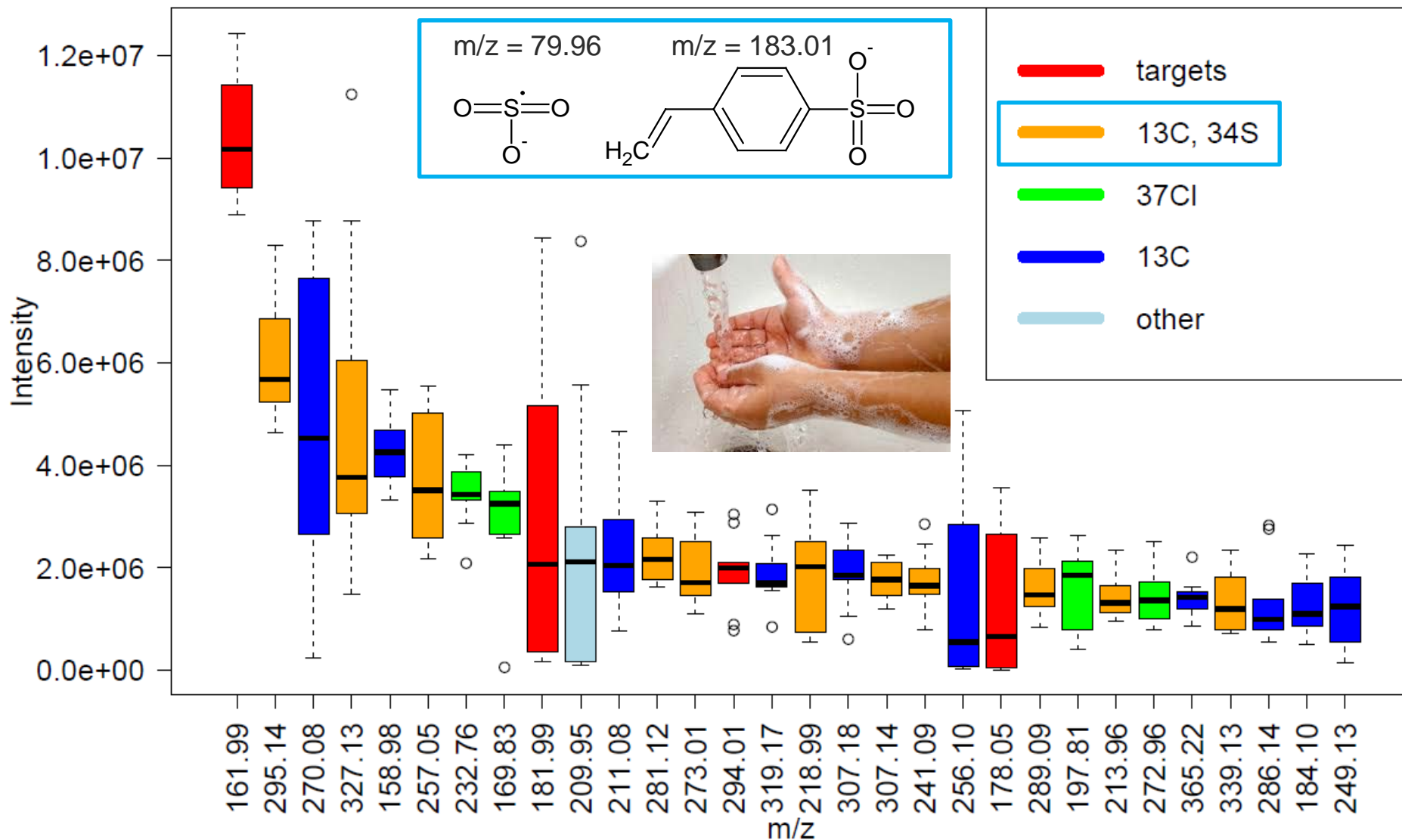
enviPat Web **nontarget**

■ Iso & Add
■ Add Only
■ Iso Only
■ Cmpts



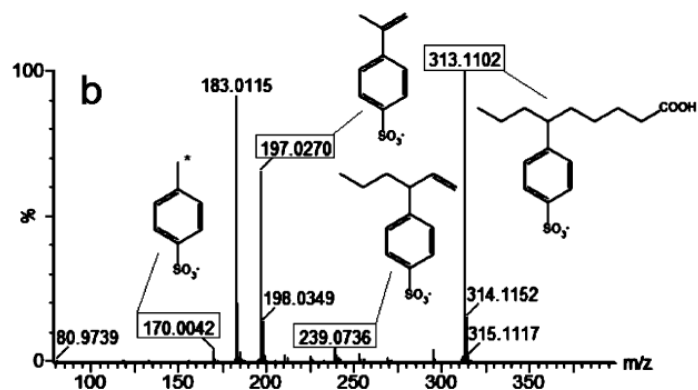
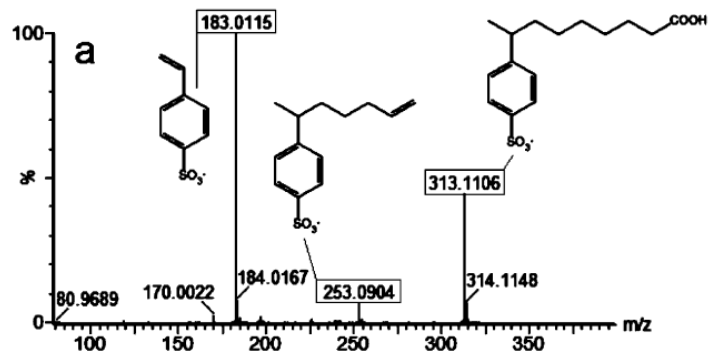


Swiss Wastewater: Top 30 Peaks (ESI-)

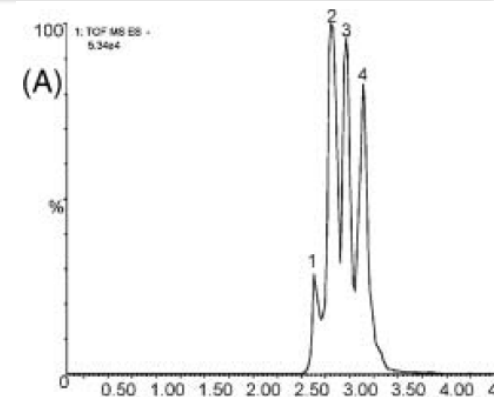


Literature sources

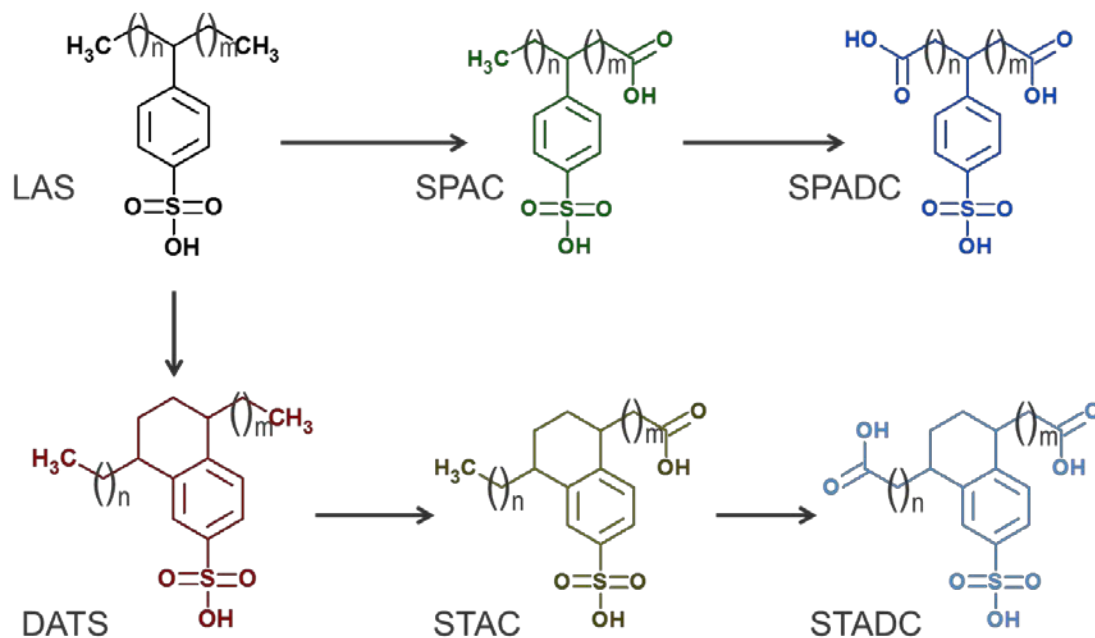
- Formulas, masses (ions), retention times and intensities
- Spectra of selected compounds (different instruments)



Lara-Martin et al. EST. 2010, 44: 1670-1676



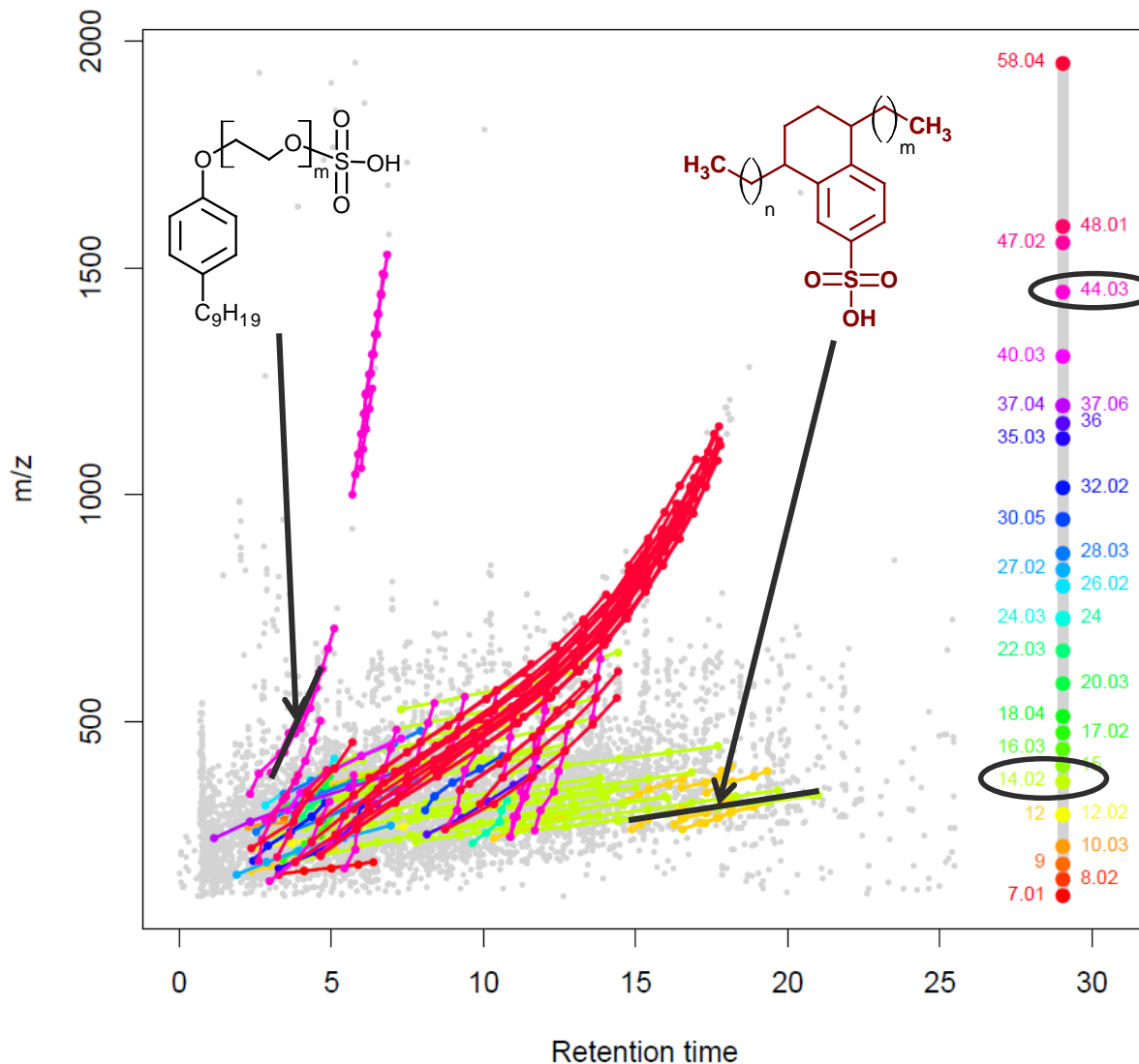
Gonzalez et al. Rapid Comm. Mass Spec. 2008, 22: 1445-54



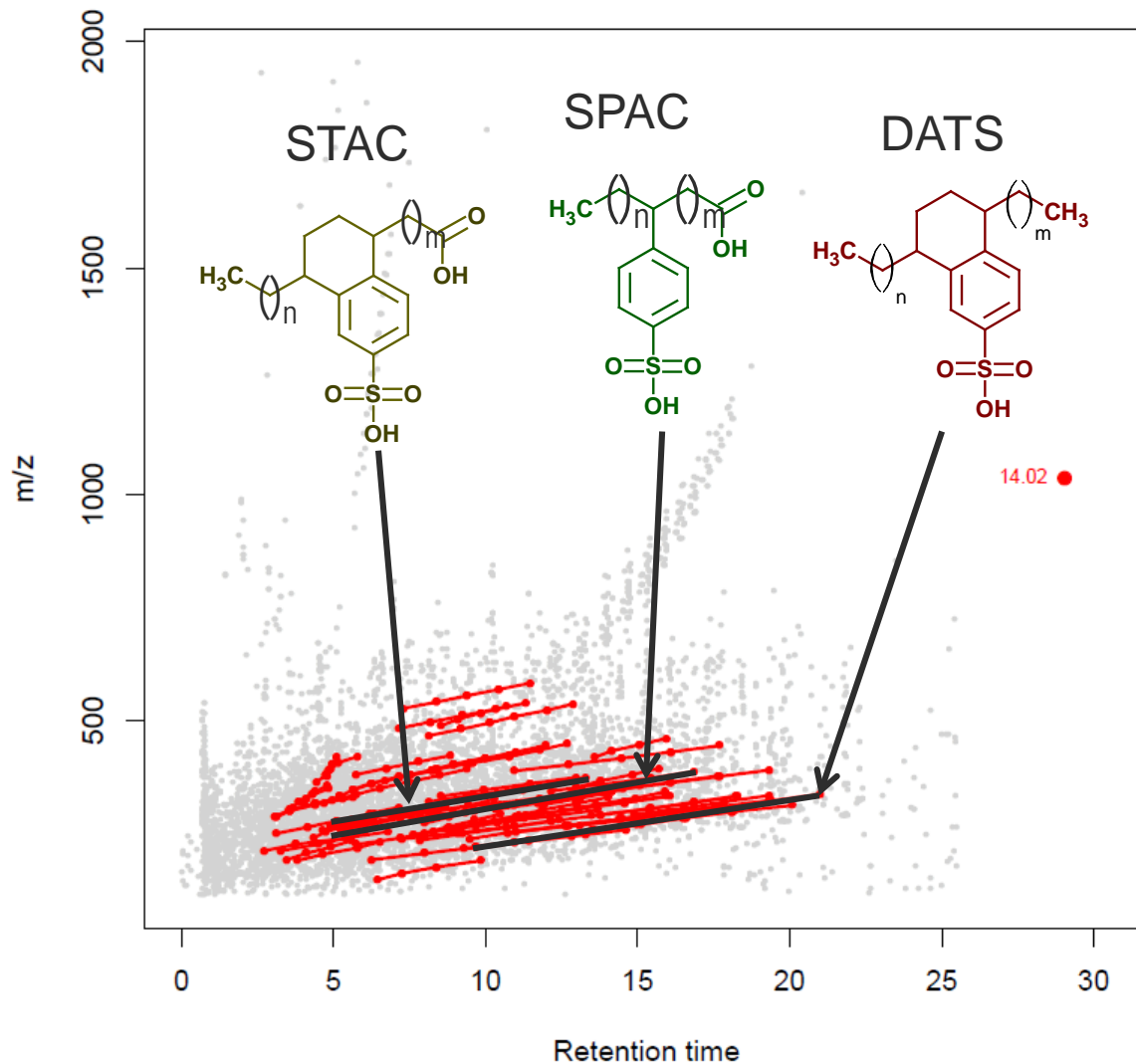
Homologous Series Detection



Search for
discrete
mass
differences



<http://www.envihomolog.eawag.ch/>

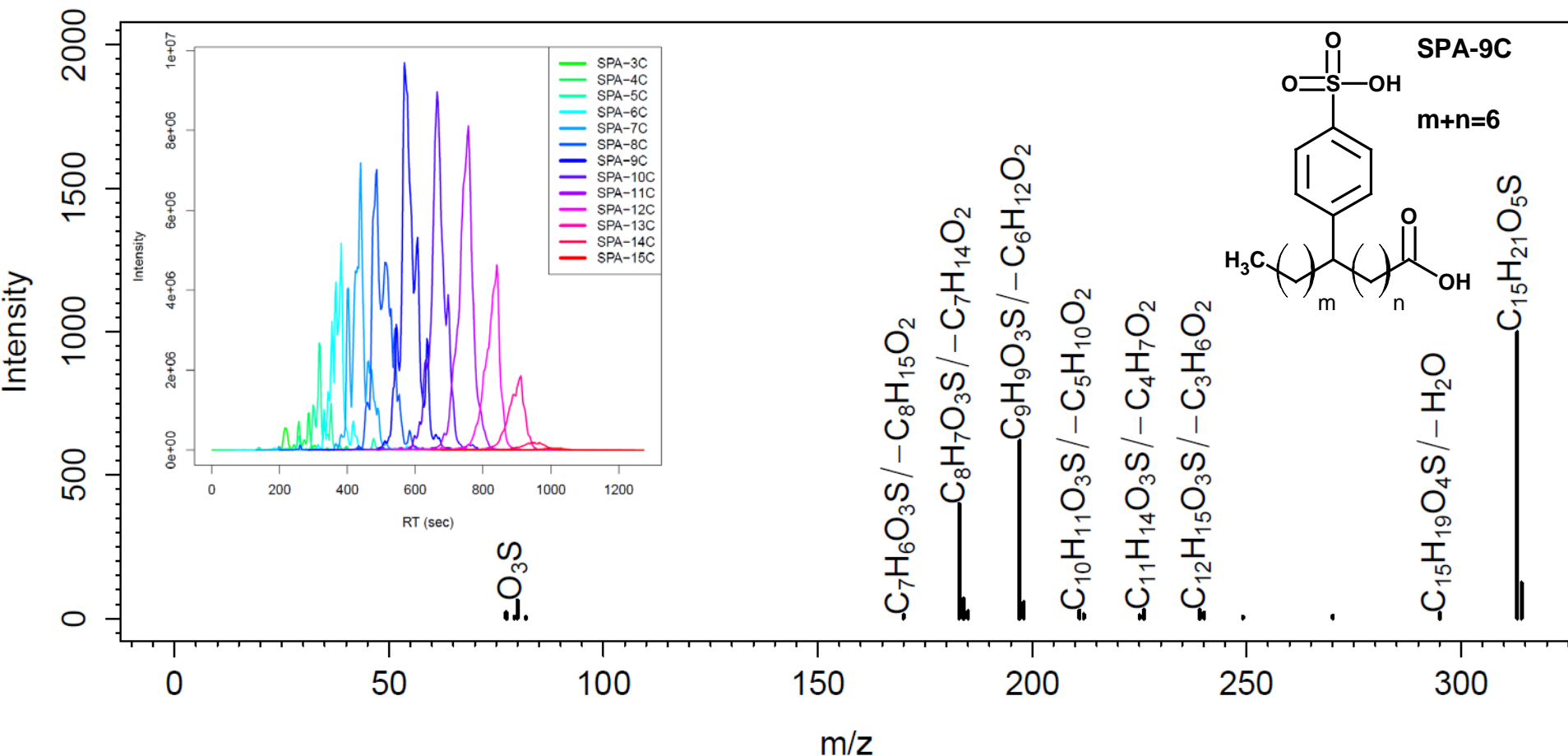


Supporting Evidence for Homologues

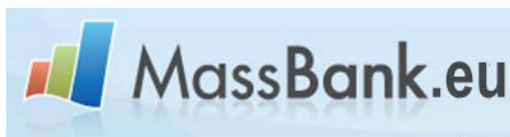
Chromatography and MS/MS Annotation

<https://github.com/MassBank/RMassBank/>

RMassBank

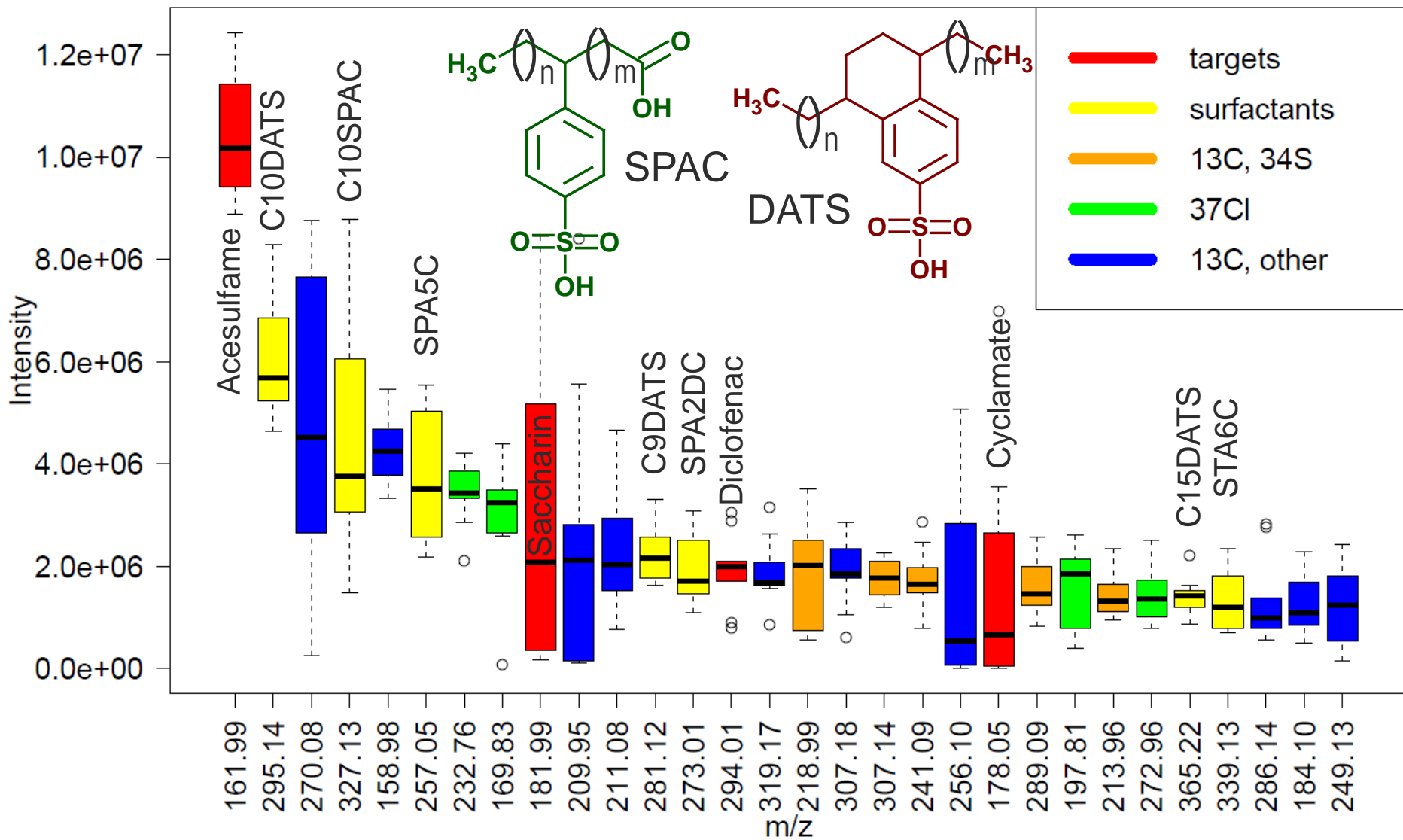


Formulas: <http://sourceforge.net/projects/genform/>
Meringer *et al.*, 2011, *MATCH* 65, 259-290
Data: Schymanski *et al.*, 2014, *ES&T*, 48:
1811-1818. DOI: 10.1021/es4044374

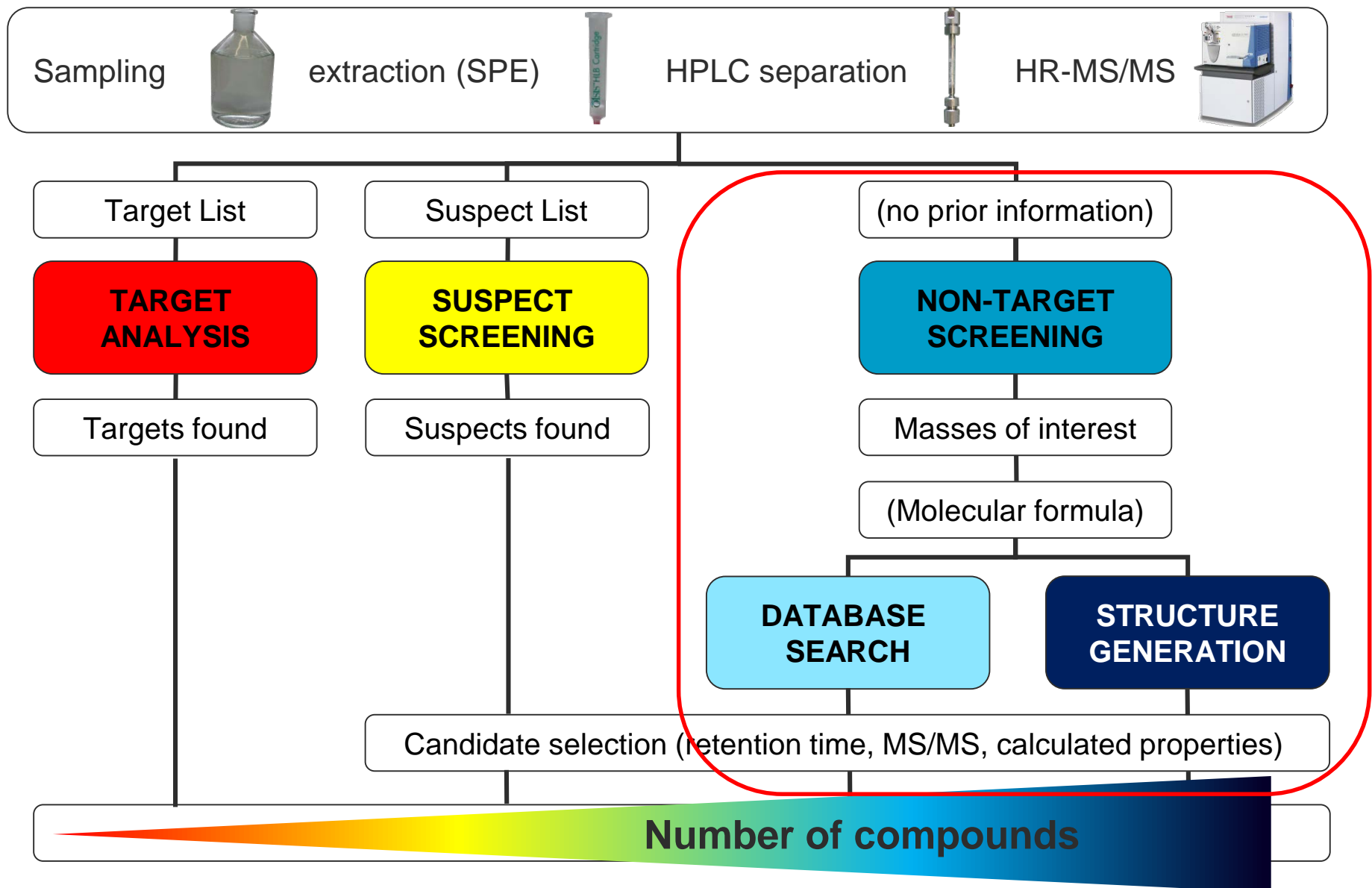


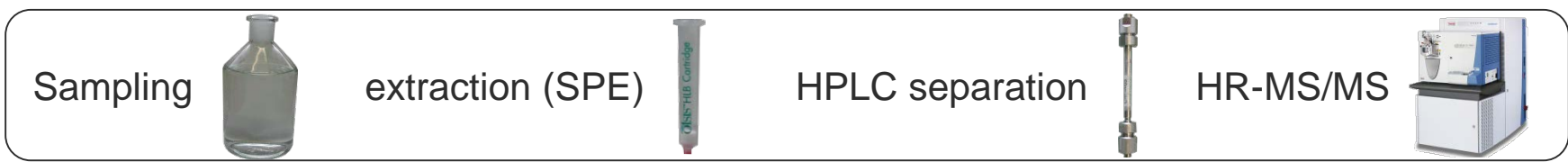
Literature: LIT00034,35
Sample: ETS00002
Standard: ETS00016,17,19,20

Swiss Wastewater: Top 30 Peaks (ESI-)



What about Non-Target Screening?





Conversion (Proteowizard) and Peak Picking (enviPick, xcms, MZmine, ...)

Detection of blank/blind/noise/internal standards; time trend analysis (enviMass)

Target List

Suspect List
(e.g. NORMAN,
LMC, Eawag-PPS,
ReSOLUTION)

NON-TARGET
SCREENING

TARGET
ANALYSIS

SUSPECT
SCREENING

Componentization
(nontarget)

Prioritization
(enviMass)

(enviMass,
vendor software)

Molecular formula
determination
(enviPat, GenForm)

Masses of interest

MS/MS Extraction
(RMassBank)



Non-target identification
(MetFrag2.3, ReSOLUTION)

Interpretation, confirmation, peak inventory, confidence and reporting

m/z $[M-H]^-$

213.9637

± 5 ppm

5 ppm

0.001 Da

ChemSpider
Search and share chemistry

or

PubChem | OPEN
CHEMISTRY
DATABASE

MetFrag

MS/MS

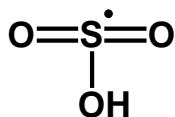
134.0054	339689
150.0001	77271
213.9607	632466

2010 => 2016

m/z [M-H]⁻
213.9637
± 5 ppm

Elements: C, N, S

5 ppm
0.001 Da



RT: 4.54 min
355 InChI/RTs

ChemSpider
Search and share chemistry

or

PubChem OPEN CHEMISTRY DATABASE

MetFrag

References
External Refs
Data Sources
RSC Count
PubMed Count



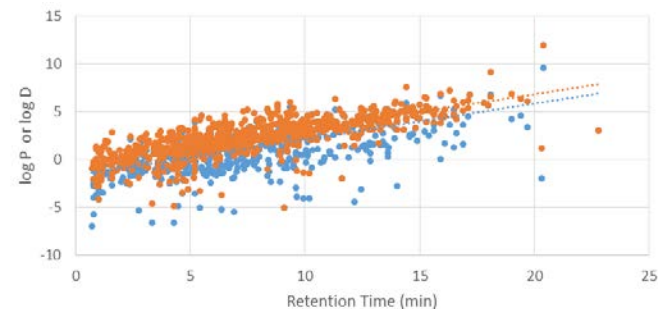
Suspect Lists
?TOFF IDENT
Chemistry Dashboard

MoNA
MassBank of North America

MassBank.eu

MS/MS

134.0054	339689
150.0001	77271
213.9607	632466



Test set of 473 Eawag Target Substances

	MetFrag 2010	MetFrag2.3 Fragments only	MetFrag2.3 +References +Retention time
ChemSpider¹			
Top 1 Ranks	73	105	420
% Top 1 Ranks	15 %	22 %	89 %
PubChem²			
Top 1 Ranks	-	30	336
% Top 1 Ranks	-	6 %	71 %

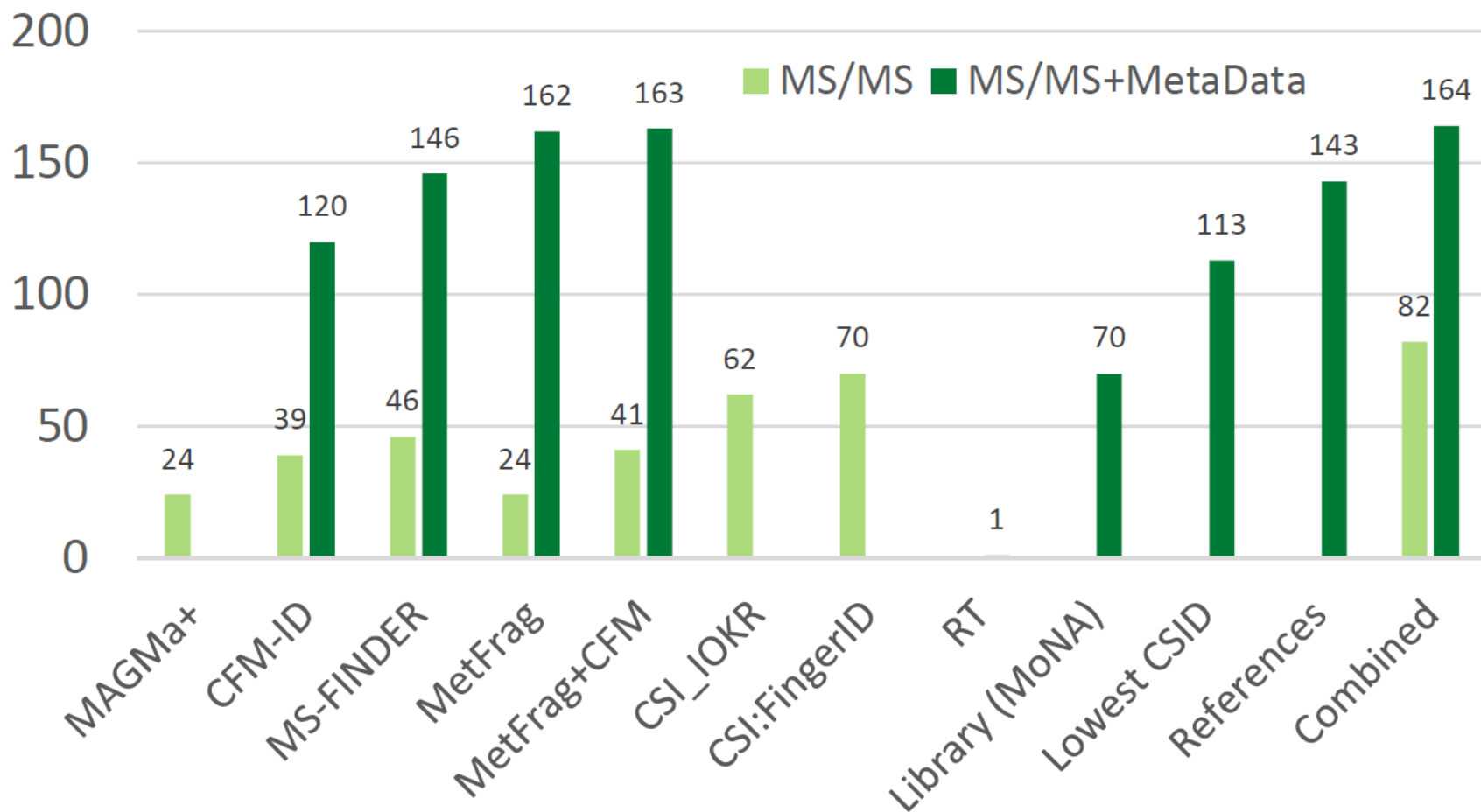
Similar results with 3 independent datasets of 310, 289 and 225 substances from Eawag and UFZ (www.massbank.eu)

¹www.chemspider.com; ~34 million entries

²<https://pubchem.ncbi.nlm.nih.gov/>; ~74 million entries

State of the Art in Small Molecule Identification

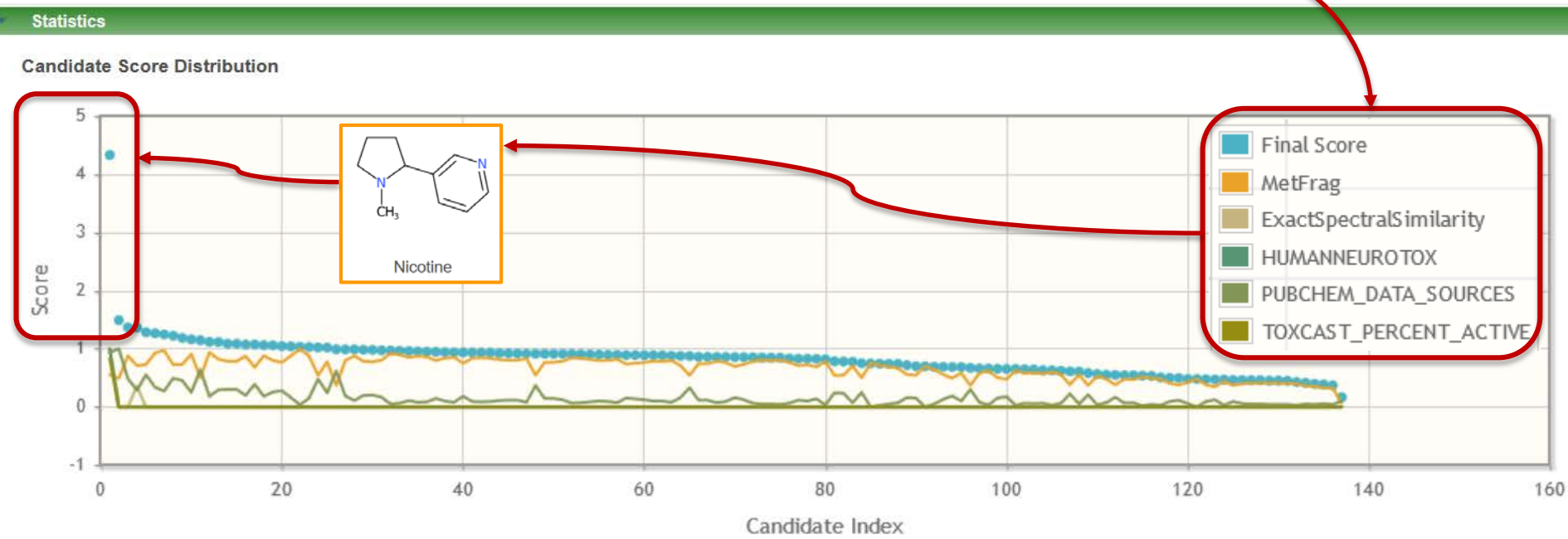
Metadata is critical to improving annotation of known unknowns!



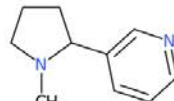
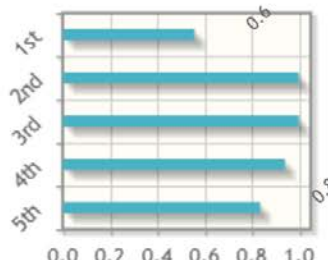
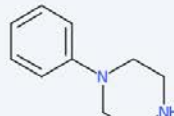
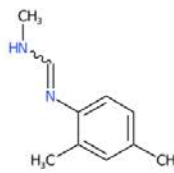
Connecting Resources in MetFrag



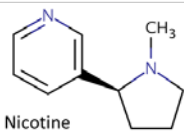
Combined evidence clearly highlights **potential neurotoxicant** among chemical candidates

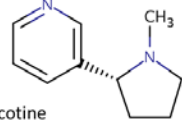


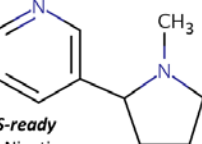
"MS-ready" Form for MetaData in MetFrag

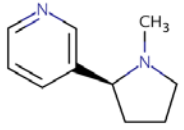
#	Molecule	Identifier	Mass	Formula	Normalized Scores	FinalScore	Details
1	 Nicotine	DTXSID1020930 DTXSID8021725 DTXSID3048154 DTXSID0046351 DTXSID6020931 DTXSID00657553 DTXSID5075319 InChIKeyBlock1 = SNICXCGAKADSCV	162.11576	C ₁₀ H ₁₄ N ₂		4.3349	Peaks: 18 / 23 Fragments Scores Download
2	 Phenylpiperazine	DTXSID40176612 DTXSID40193102 DTXSID90216632 DTXSID50291046 DTXSID00293111 DTXSID50296613 InChIKeyBlock1 = YZTJYBJCZXZGCT	162.11576	C ₁₂ H ₁₆ N ₂			
3	 N-(2,4-Dimethylphenyl)-N-methylformamide	DTXSID1037696 DTXSID10199510 InChIKeyBlock1 = JIIOLEGNERQDIP	162.11576	C ₁₂ H ₁₆ N ₂			

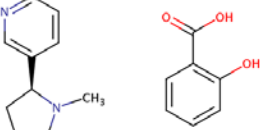
LEGEND: Name, SMILES
DTXSID | InChIKey 1st Block
CAS | Monoiso. Mass | logP | Sources
Data on: Toxicity | Exposure | Bioassays

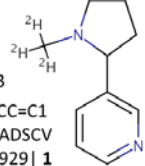

Nicotine
CN1CCC[C@H]1C1=CN=CC=C1
DTXSID1020930 | SNICXCGAKADSCV
54-11-5 | **162.1157** | 0.929 | **72**
Tox: **yes** | Expo: **yes** | Bioassay: **yes**


D-Nicotine
CN1CCC[C@@H]1C1=CN=CC=C1
DTXSID0046351 | SNICXCGAKADSCV
25162-00-9 | **162.1157** | 0.929 | **20**
Tox: **no** | Expo: **yes** | Bioassay: **yes**


MS-ready
DL-Nicotine
CN1CCCC1C1=CN=CC=C1
DTXSID3048154 | SNICXCGAKADSCV
22083-74-5 | **162.1157** | 0.953 | **9**
Tox: **yes** | Expo: **no** | Bioassay: **yes**


HCl
Nicotine hydrochloride
Cl.CN1CCC[C@H]1C1=CN=CC=C1
DTXSID6020931 | HDJBTCAJIMNXEW
2820-51-1 | **198.0924** | 0.929 | **9**
Tox: **no** | Expo: **yes** | Bioassay: **yes**


Benzoic acid, 2-hydroxy-, compd. with
3-[(2S)-1-methyl-2-pyrrolidinyl]pyridine (1:1)
OC(=O)C1=C(O)C=CC=C1.CN1CCC[C@H]1C1=CN=CC=C1
DTXSID5075319 | AIBWPBUAKCMKNS
29790-52-1 | **300.1474** | 0.929 | **6**
Tox: **no** | Expo: **yes** | Bioassay: **no**

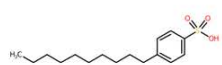
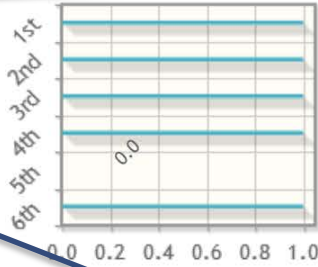
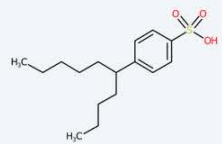
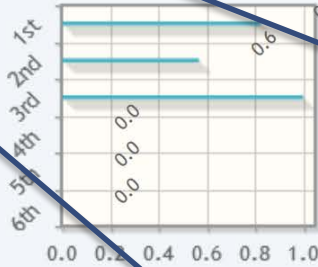
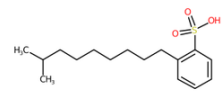
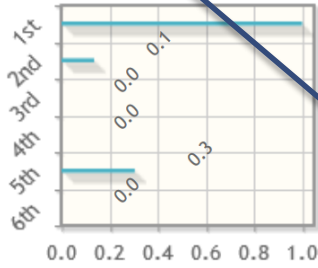

DL-Nicotine-d3
[2H]C([2H])([2H])N1CCCC1C1=CN=CC=C1
DTXSID80442666 | SNICXCGAKADSCV
69980-24-1 | **165.1345** | 0.929 | **1**
Tox: **no** | Expo: **no** | Bioassay: **no**

"MS-ready" Form for MetaData in MetFrag



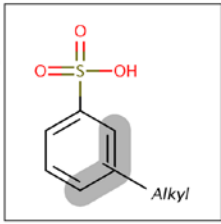
MetFrag <https://msbi.ipb-halle.de/MetFragMSready/> Search

Results

#	Molecule	Identifier	Mass	Formula	Normalized Scores
1	 (C10-C16) Alkylbenzenesulfonic acid	DTXSID2028723 DTXSID7059696 DTXSID5041647 InChIKeyBlock1 = UASQKKHYUPBQJR	298.16027	C ₁₆ H ₂₆ O ₃ S	
2	 Alkylbenzenesulfonate, linear	DTXSID3020041 DTXSID70881146 InChIKeyBlock1 = KIIODTIGNUGDLO	298.16027	C ₁₆ H ₂₆ O ₃ S	
7	 Benzenesulfonic acid, isodecyl-, compd. with 2,2'-iminobis[ethanol] (1:1)	DTXSID2070762 InChIKeyBlock1 = CKNPXKLNOLAXDF	298.16027	C ₁₆ H ₂₆ O ₃ S	

(C10-C16) Alkylbenzenesulfonic acid
68584-22-5 | DTXSID2028723

Search by DSSTox_Substance_id Found 1 result for DTXSID2028723

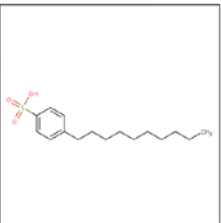


Intrinsic
Molec
Avera
Mono

Structur
Presenc
Record
Quality

4-Decylbenzenesulfonic acid
140-60-3 | DTXSID7059696

Search by DSSTox_Substance_id Found 1 result for DTXSID7059696



Download

Items: 4 / 5

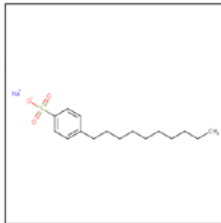
ments

pres

Download

Sodium 4-decylbenzenesulfonate
2627-06-7 | DTXSID5041647

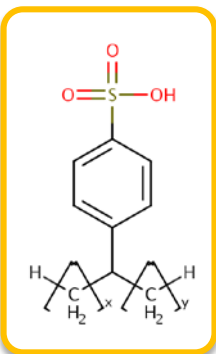
Search by DSSTox_Substance_id Found 1 result for DTXSID5041647



Connecting and Enhancing Open Resources

- Sharing knowledge is a **win-win situation**

2014



MassBank.eu

norman
suspects

EPA
Chemicals

2015: found in waters across Europe



2016: 1 datapoint cross-annotates 3072 in GNPS

Hits in GNPS Massive datasets:

Surfactants: <http://goo.gl/7sY9Pf>

2017: Early-Warning System is born

normanews

2018: Highlighted in



CHEMISTRY
Early warning about emerging contaminants



“Live” retrospective screening of known and unknown chemicals in European samples (various matrices)



www.norman-data.eu

NORMAN Digital Sample Freezing Platform

Main Page

Batch mode

Contributed Samples

Results

Chromatograms

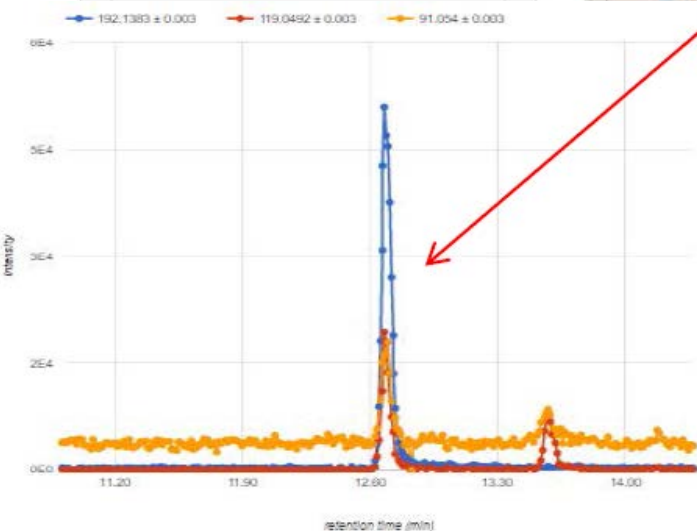
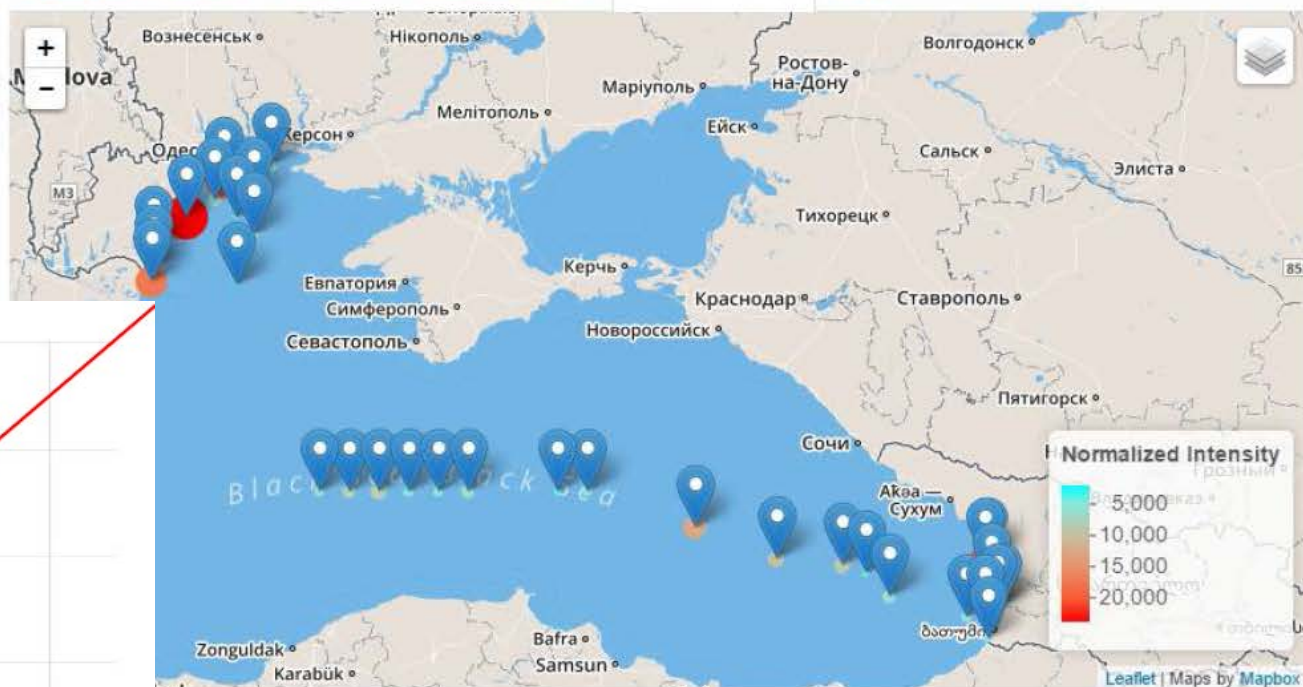
Interactive Map

Help

Choose Emerging Substance or input mass of interest and experimental RTI

Substance name or CAS or StdInChIKey

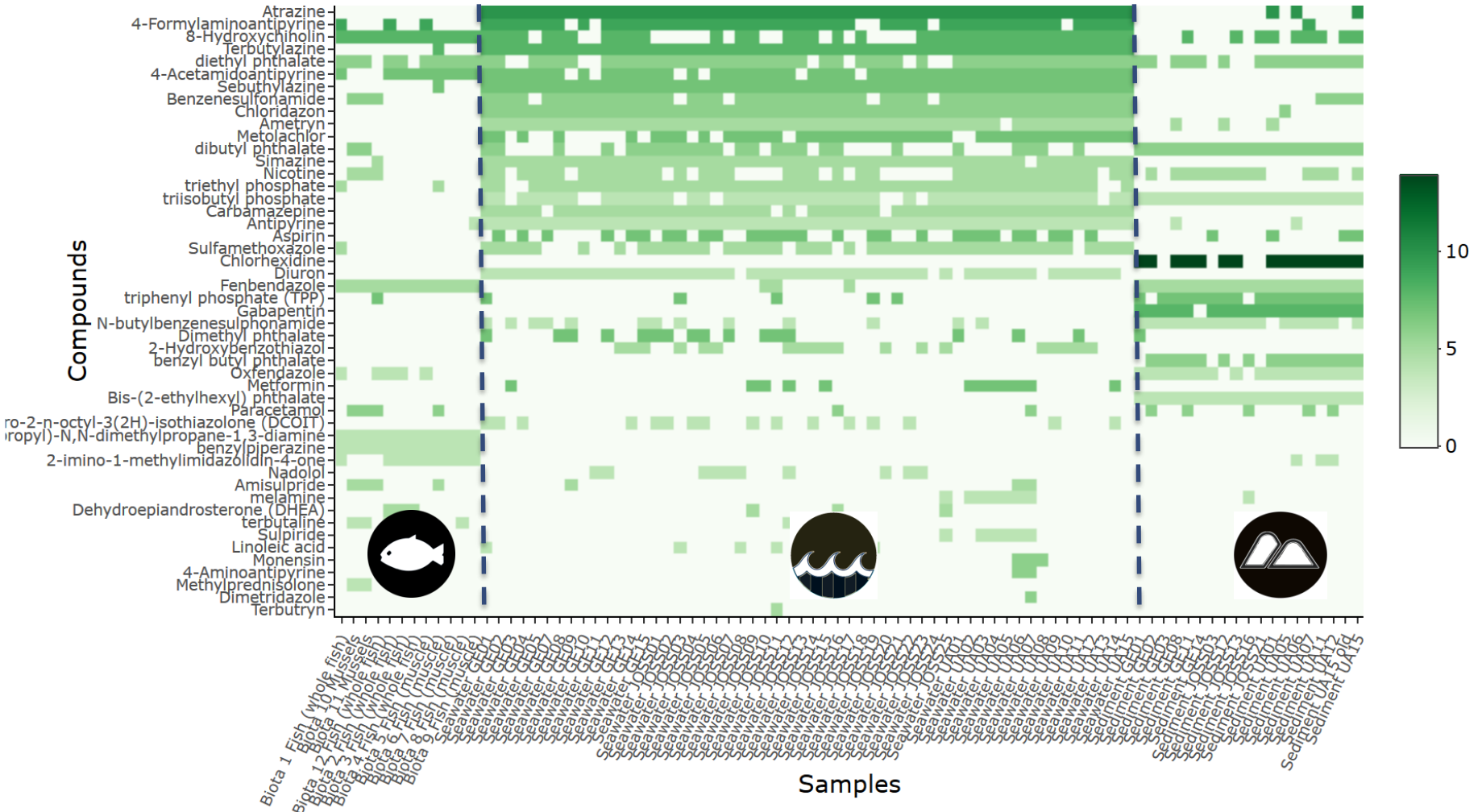
DEET [134-62-3]
[MMOXZBCLCQITDF-
UHFFFAOYSA-N]



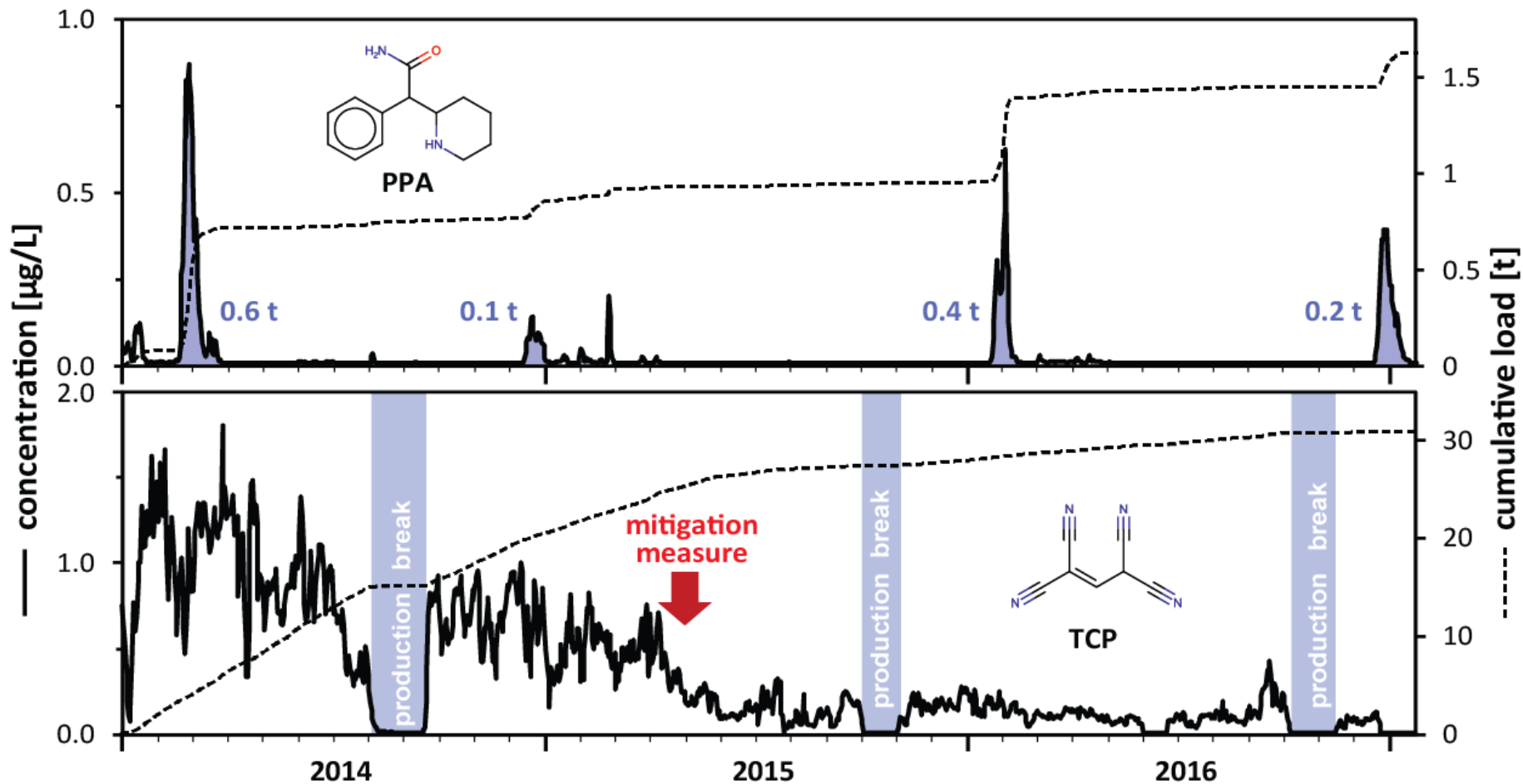
Retrospective screening of REACH chemicals in Black Sea samples (various matrices)

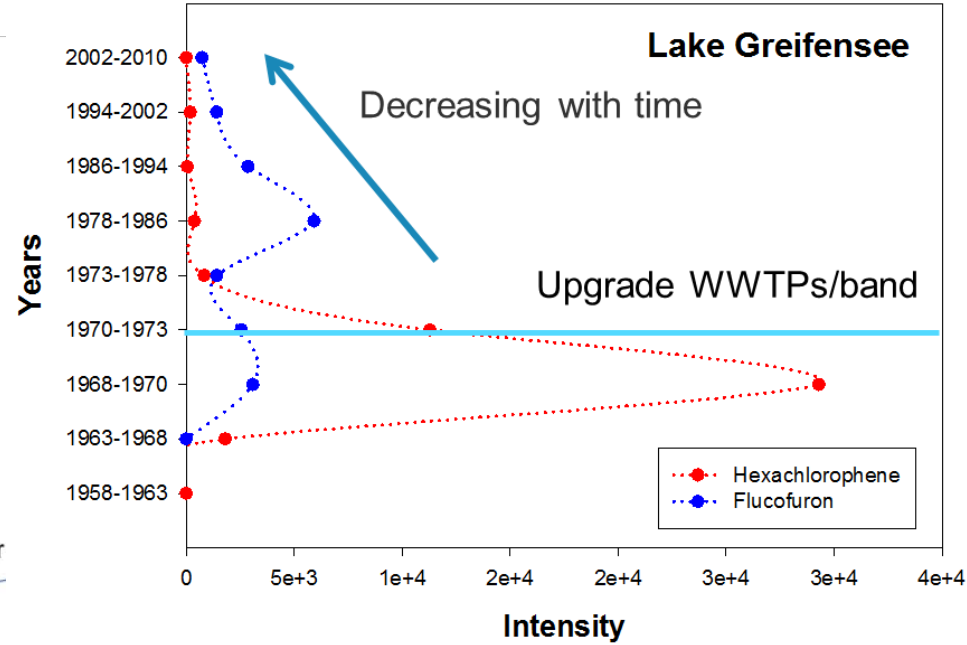
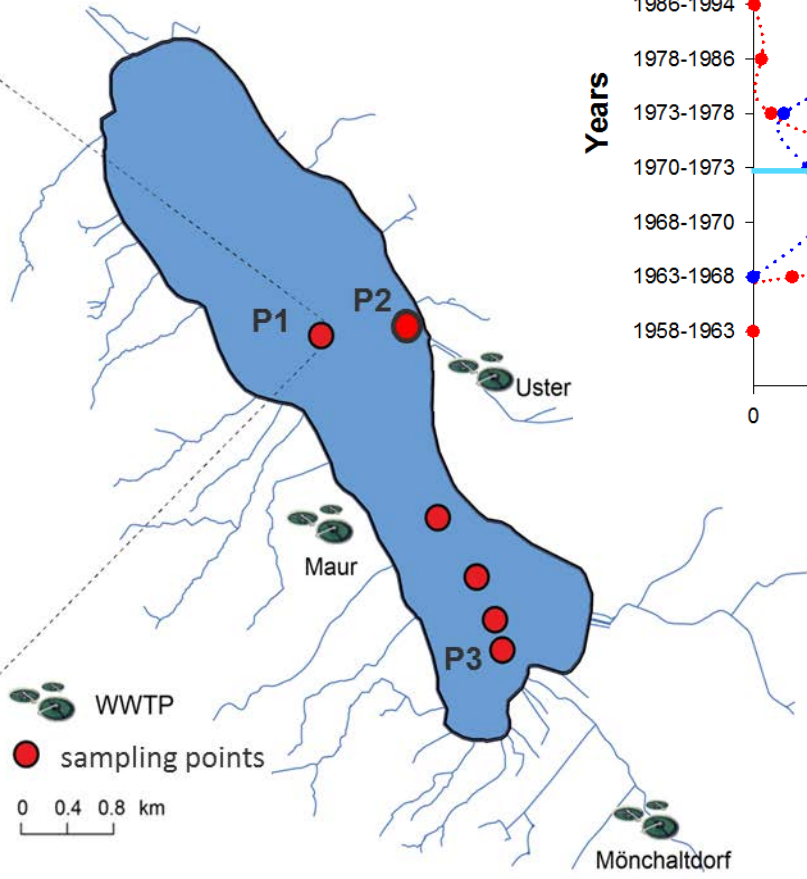
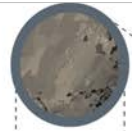


Occurrence Results

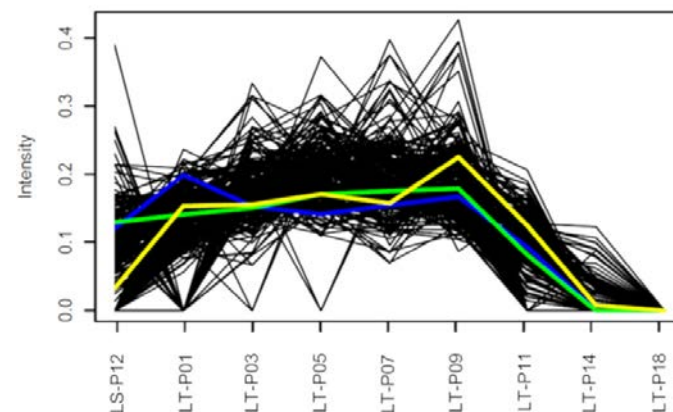
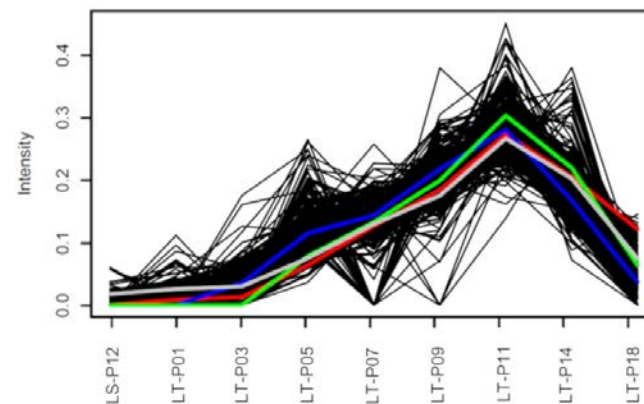
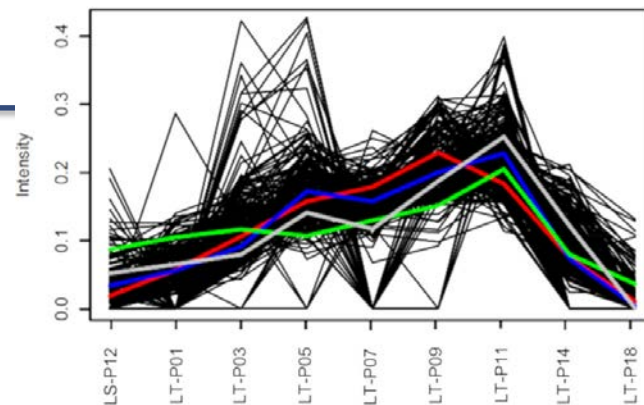
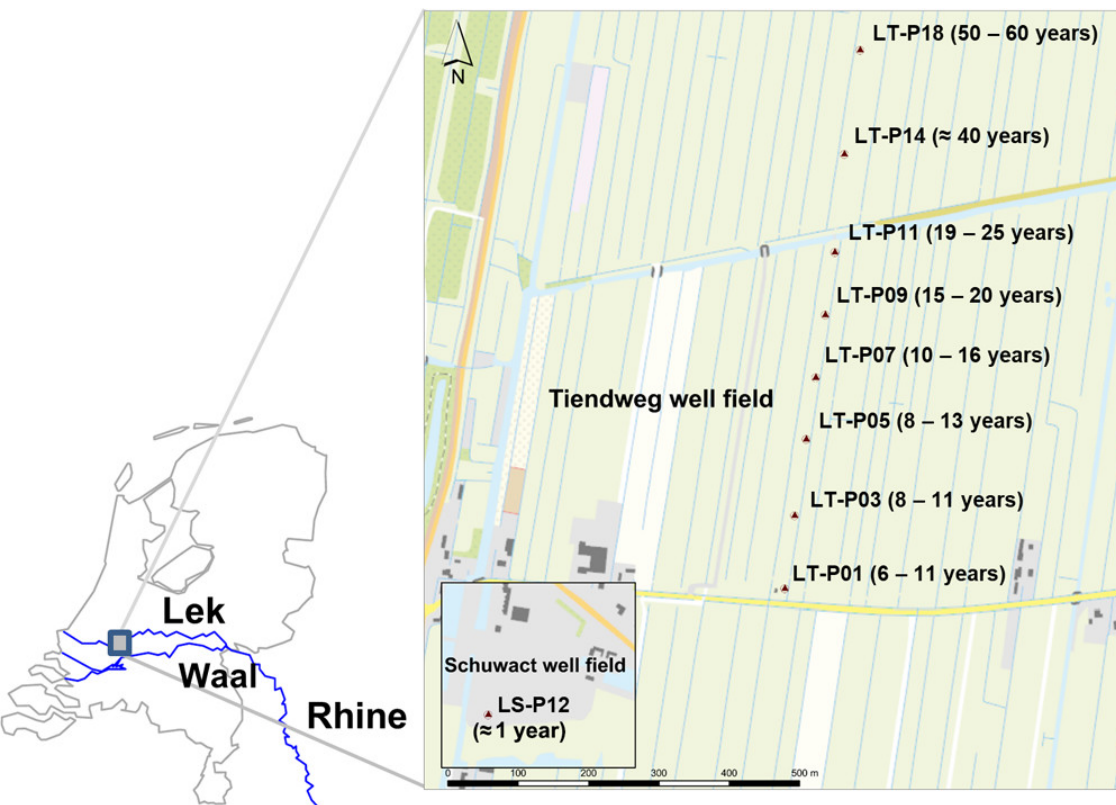


Previously unknown chemicals detected due to “stand-out” patterns





Micropollutant Time Trends in Riverbank Filtration Systems



New MetaData: Disease-Specific Reference Counts

https://comptox.epa.gov/dashboard/chemical_lists/litminedneuro

Chemical	CAS RN	DSSToxID	PMID Ct	Seizures	Nervous System Diseases	Peripheral Nervous System Diseases	Brain Diseases	Muscular Diseases	Basal Ganglia Diseases	Parkinson Disease, Secondary	Coma	Hallucinations	Tremor	Memory Disorders	Central Nervous
Cisplatin	15663-27-1	DTXSID4024983	1032	20	47	140	13	0	4	1	1	0	1	2	4
Ethanol	64-17-5	DTXSID9020584	768	100	23	11	18	26	1	3	20	6	17	54	2
Lead	7439-92-1	DTXSID2024161	740	28	107	68	102	4	2	2	1	3	4	19	30
Lithium	7439-93-2	DTXSID5036761	689	30	50	9	22	5	36	13	25	6	93	12	15
Valproic Acid	76584-70-8	DTXSID70227388	666	32	10	3	65	6	10	18	45	5	18	4	2
1-Methyl-4-phen	28289-54-5	DTXSID8040933	638	1	24	0	11	0	6	289	0	0	5	0	1
Vincristine	2068-78-2	DTXSID8044331	567	17	59	125	15	5	1	1	5	3	2	1	8
Phenytoin	57-41-0	DTXSID8020541	560	37	24	25	16	9	3	1	9	3	8	4	6
Haloperidol	52-86-8	DTXSID4034150	555	6	6	1	10	6	153	51	4	4	11	1	0
Cocaine	50-36-2	DTXSID2038443	530	151	16	0	8	0	2	3	3	8	6	12	11
Aspirin	50-78-2	DTXSID5020108	489	8	3	0	3	2	2	0	9	4	1	0	5
Paclitaxel	33069-62-4	DTXSID9023413	485	4	43	217	9	14	0	0	0	0	0	1	2
Aluminum	7429-90-5	DTXSID3040273	477	13	41	1	105	4	0	0	1	0	1	13	12
Lidocaine	6108-05-0	DTXSID80209953	464	150	26	15	3	2	0	0	8	4	6	2	10
Methotrexate	59-05-2	DTXSID4020822	451	17	25	1	79	4	0	1	5	0	1	9	18
Mercury	7439-97-6	DTXSID1024172	450	6	79	22	23	2	3	5	2	2	38	7	25

New MetaData: Disease-Specific Reference Counts


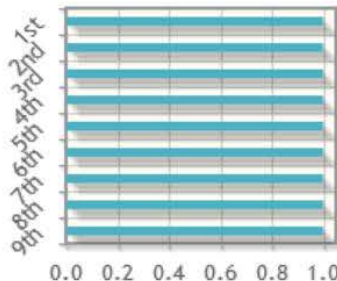
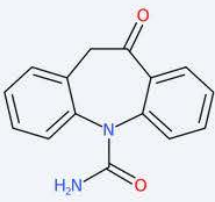
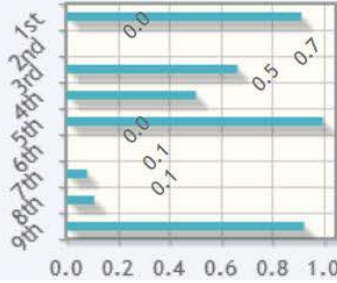
<https://msbi.ipb-halle.de/MetFrag/>

Weights

MetFrag (1st)	<input type="range" value="100"/>	100 %
ExactSpectralSimilarity (2nd)	<input type="range" value="100"/>	100 %
Basal Ganglia Diseases (3rd)	<input type="range" value="100"/>	100 %
DATA_SOURCES (4th)	<input type="range" value="100"/>	100 %
Parkinson Disease, Secondary (5th)	<input type="range" value="100"/>	100 %
Peripheral Nervous System Diseases (6th)	<input type="range" value="100"/>	100 %

Download Results

Filter Candidates by explained MS/MS Peaks

#	Molecule	Identifier	Mass	Formula	Normalized Scores	FinalScore	Details
1	 <p>Phenytoin</p>	<p><u>DTXSID8020541</u></p> <p>InChIKeyBlock1 = <u>CXOFVDLJLONNDW</u></p>	252.08992	C ₁₅ H ₁₂ N ₂ O ₂		9.0	<p>Peaks: 9 / 14</p> <p><input type="button" value="Fragments"/></p> <p><input type="button" value="Scores"/></p> <p><input type="button" value="Download"/></p>
2	 <p>oxcarbazepine</p>	<p><u>DTXSID0045703</u></p> <p>InChIKeyBlock1 = <u>CTRLABGOLIVAIV</u></p>	252.08992	C ₁₅ H ₁₂ N ₂ O ₂		4.2006	<p>Peaks: 11 / 14</p> <p><input type="button" value="Fragments"/></p> <p><input type="button" value="Scores"/></p> <p><input type="button" value="Download"/></p>

*Future work: use results of **unknowns** to drive prioritization efforts*

EMBLAS-II: EU/UNDP project Improving Environmental Monitoring in the Black Sea Non-target Screening Data

NTS Data

Joint Black Sea Surveys (JBSS)
year

2017

Choose compound

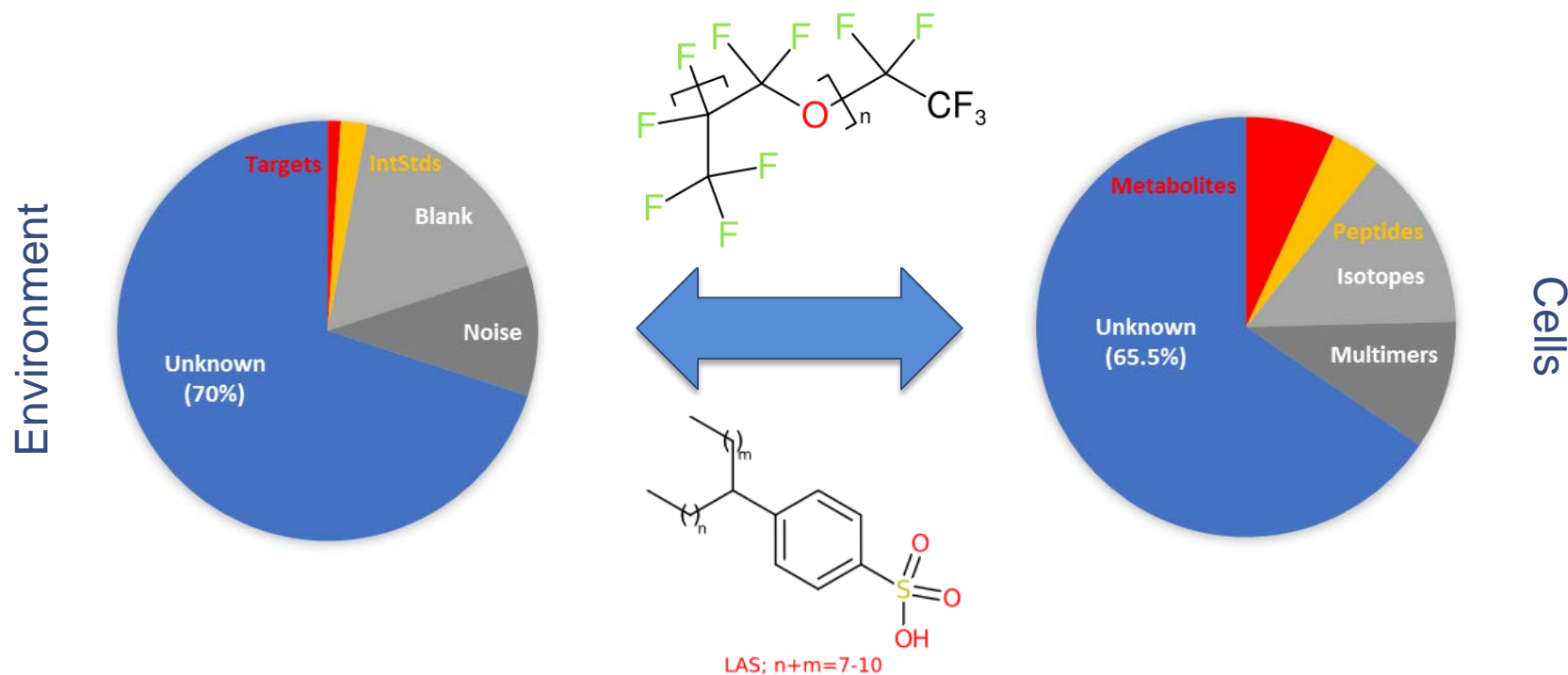
Positive ESI mz=354.2876
rt=10.87



Take Home Messages

Unknowns and High Resolution Mass Spectrometry

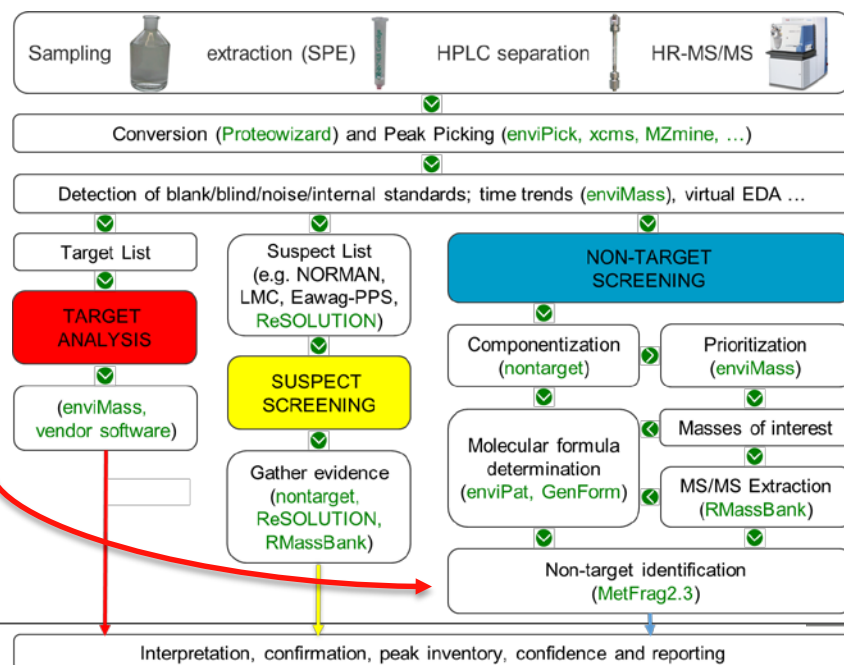
- Over 60 % of HR-MS peaks are potentially *relevant* but *unknown*



Take Home Messages

Unknowns and High Resolution Mass Spectrometry

- Over 60 % of HR-MS peaks are potentially **relevant** but **unknown**
- **Annotating unknowns** requires data and evidence from **many different sources**
 - Many excellent workflows available to collate this information
 - Incorporation of **all available metadata** is critical to success!
 - *E.g.* MetFrag has greatly improved the speed and success of tentative identification of “known unknowns”: **15 % => 89 % Ranked Number 1**
- <https://ipb-halle.github.io/MetFrag/>



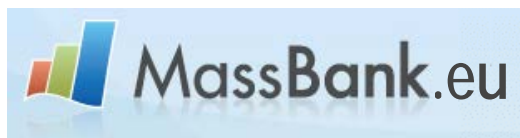
Take Home Messages

Unknowns and High Resolution Mass Spectrometry

- Over 60 % of HR-MS peaks are potentially **relevant** but **unknown**
- **Annotating unknowns** requires data and evidence from **many different sources**
- Exchange expert knowledge worldwide
 - Community efforts contribute greatly to improved cross-annotation
 - Information in the public domain helps everyone!
 - You never know when it will help you 😊



Acknowledgements



solutions



emma.schymanski@uni.lu

@ESchymanski

Further Information:

<https://massbank.eu/MassBank/>

<https://ipb-halle.github.io/MetFrag/>

<https://www.norman-network.com/nds/SLE/>


[https://wwwen.uni.lu/lcsb/research/
environmental_cheminformatics](https://wwwen.uni.lu/lcsb/research/environmental_cheminformatics)





Community Efforts!



 MassBank consortium

