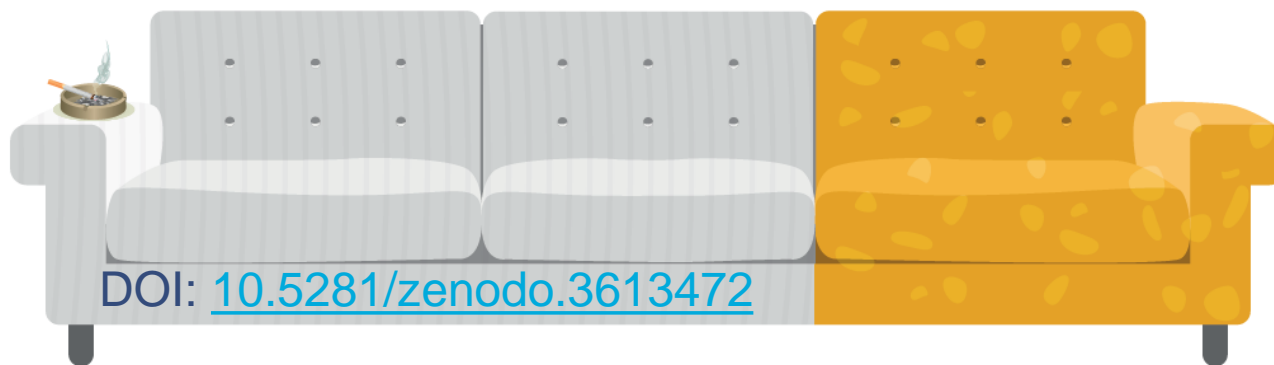


# Environmental Cheminformatics: *Case study of Thirdhand Smoke in House Dust.*



Emma L. Schymanski<sup>1</sup>, Sonia Torres<sup>2</sup>, Noelia Ramirez<sup>2</sup>

<sup>1</sup>FNR ATTRACT Fellow; Luxembourg Centre for Systems Biomedicine, University of Luxembourg

<sup>2</sup>Metabolomics Core, IISPV-University Rovira i Virgili, Tarragona, Spain






[emma.schymanski@uni.lu](mailto:emma.schymanski@uni.lu)


Twitter: [@ESchymanski](https://twitter.com/ESchymanski) [@soniatorres](https://twitter.com/soniatorres) [@noeliaramz](https://twitter.com/noeliaramz)


...plus many other colleagues who have contributed over the years!

## LET'S MAKE IT HAPPEN


### Members with access to **Environmental Cheminformatics**

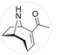
-  **Adelene Lai** @adelene.lai  
Given access 2 months ago
-  **Anjana Elapavalore** @anjana.elapavalore  
Given access 4 weeks ago
-  **Corey Griffith** @corey.griffith  
Given access 2 months ago
-  **Emma Schymanski** @emma.schymanski It's you  
Given access 2 months ago
-  **German Andres Preciat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


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

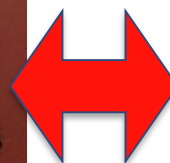
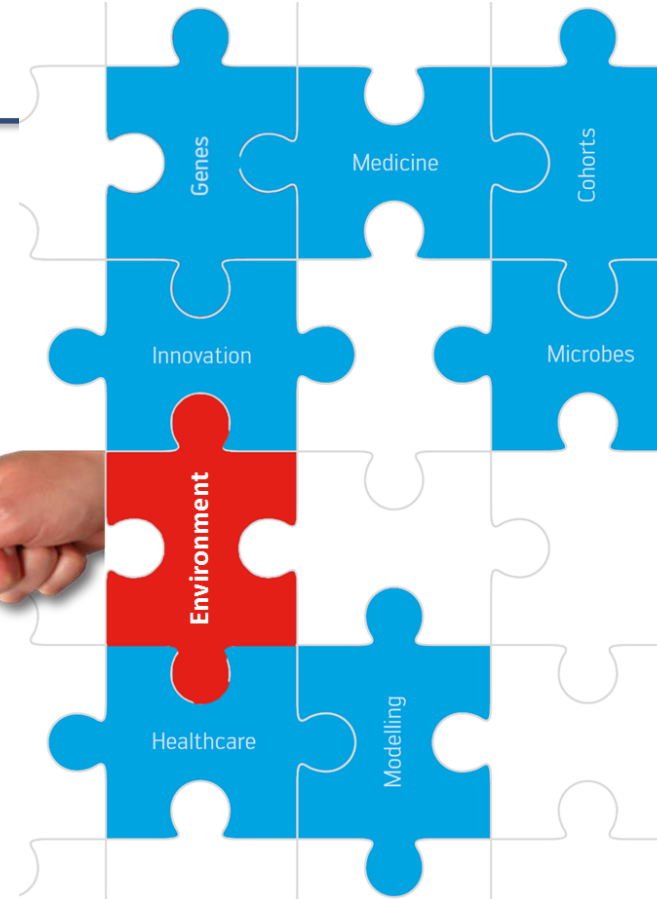
 **Jessy Krier** @jessy.krier  
Given access 1 week ago

 **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 ?

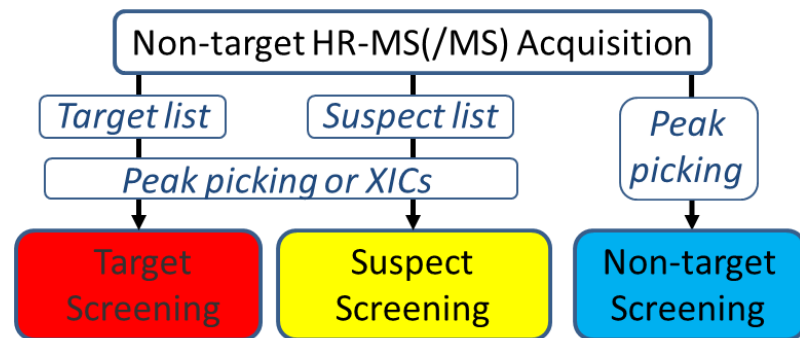


# Environmental Cheminformatics: Case Study of

## Thirdhand Smoke (THS) in House Dust

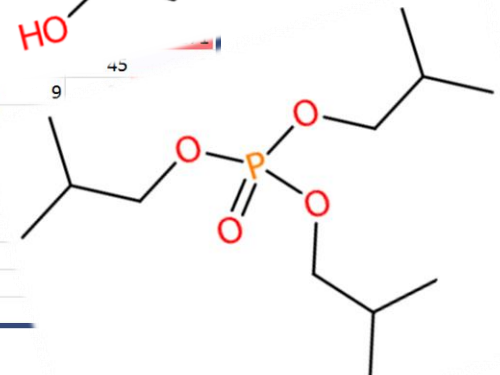
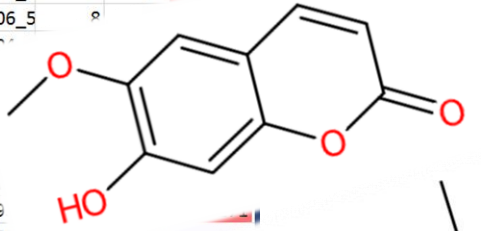
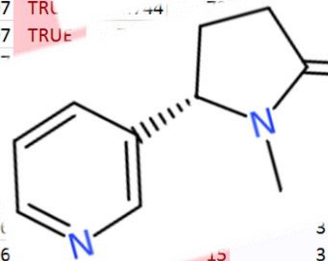
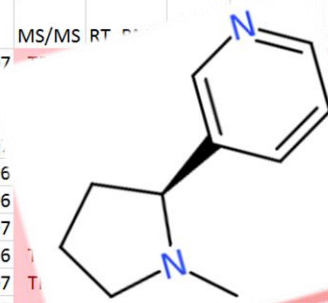


+



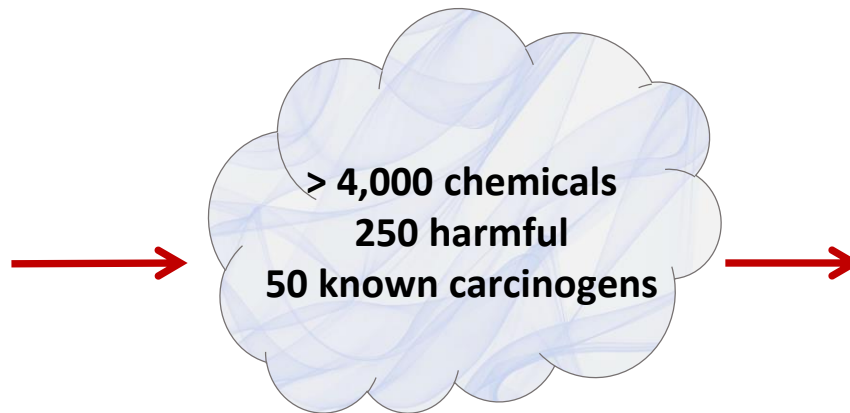
=

ID	mz	Name	RT	Int	MS/MS	RT	SMILES	Name_maxScore	ExplPeaks	#Peaks	#Peaks	max	max
3027	163.1221	FT0532	1.073	9E+07			CN1CCC	Nicotine	58.0658	9	28	0.99701	
3008							N1CC(=	Creatinine	57.0454	1	3	14	0.54
3131							C(C=CC	all-trans-Retinoic acid	57.0706	5	8		
3321							CCCC	Didecyl phthalate	69.070				
3484							C(C)(C	2,2'-Oxamidodiethyl bis[3-	57.				
3206	346.1096	FT3193	21.08	3E+06			N(CCN	Nitralin	NA				
3044							CC	N,N-Diethylnicotinamide	78.0				
3006							CC	Diethylene glycol	NA				
3043							CCN(C	N,N-Diethylnicotinamide	84.08				
3055							=C	Scopoletin	53.039				
3046	183.0796	FT0741	24.77	1E+07			1=	Benzophenone	50.015				
3039							C	Cotinine	53.0393	4	9		
3038								4-Methyl-1-phenylpyrazoli	53.0393	7			
3183								Dihexyl phthalate	54.0452				
3095								Triisobutyl phosphate	57.0707				
3020	151.1111	FT0427	13.16	2E				tert-Butylphenol	NA				
3029								4-Trimethoxybenzene	NA				
3123								Octadecyl isocyanate	NA				
3128							3 CCCCCC	Octadecanoic acid, hydraz	NA				
3172							3 CCCCCC	Stearylbenzene	NA				



# What is Thirdhand Smoke (THS)?

- **Thirdhand smoke** is the long-lasting residue resulting from second-hand smoke that accumulates in dust, in objects and on surfaces in indoor environments where tobacco has been smoked (WHO, 2017).



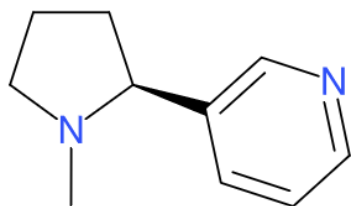
**Tobacco smoke  
(second hand smoke)**



**Thirdhand Tobacco  
Smoke (THS)**

- Poorly studied, underestimated exposure route

# What is Third Hand Smoke (THS)?



**Nicotine**

**OXIDATION** →

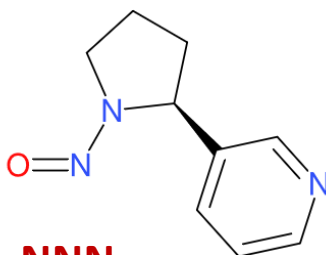
**TOBACCO-SPECIFIC  
NITROSAMINES  
(TSNAs)**



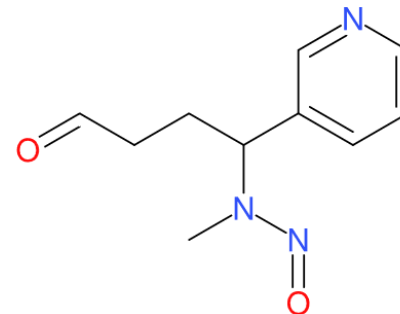
International Agency  
Research on Cancer



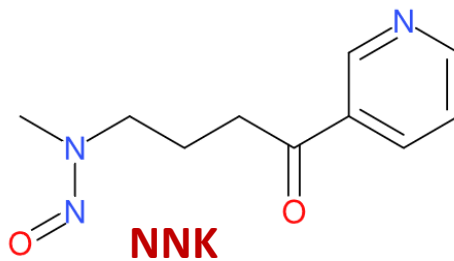
World Health  
Organization



**NNN**



**NNA  
(exclusive in THS)**



**NNK**

# Main Exposure Pathways to THS



**Dermal  
Absorption**



**Non-dietary  
Ingestion**



**Inhalation**



**Children < 5 years**

# The challenge (for NTS)? Smoke is everywhere ...

## ROME (2011)



**TPM (ng/m<sup>3</sup>)**

Nicotine 1,700-4,800

*Cecinato et al. Environ. Pollut. (2012)*

## LONDON (2012)



**PM<sub>2.5</sub> (ng/m<sup>3</sup>)**

Nicotine 21 (max. 118)

TSNAs 1.2 (max. 6.2)

*Farren et al. Environ. Sci. Technol. (2015)*

## TARRAGONA (2013)



**TPM (ng/m<sup>3</sup>)**

Nicotine 4 (max. 12.5)

*Aragón et al. JCA (2013)*

## SAN FRANCISCO (2016)



**TPM (ng/m<sup>3</sup>)**

Nicotine 339 (13% RSD)

TSNAs 0.8 (20% RSD)

*Peyton et al. Chem. Res. Toxicol. (2017)*

# THS in Tarragona, Spain



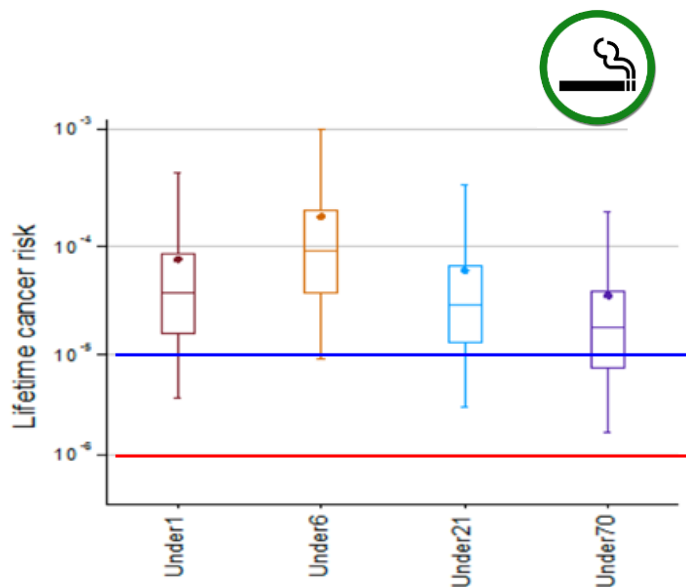
**SMOKER'S HOMES**  
(n=22)

Nicotine 26  $\mu\text{g/g}$   
TSNAs 0.5  $\mu\text{g/g}$



**NON-SMOKER'S HOMES**  
(n=24)

Nicotine 3.3  $\mu\text{g/g}$   
TSNAs 0.09  $\mu\text{g/g}$

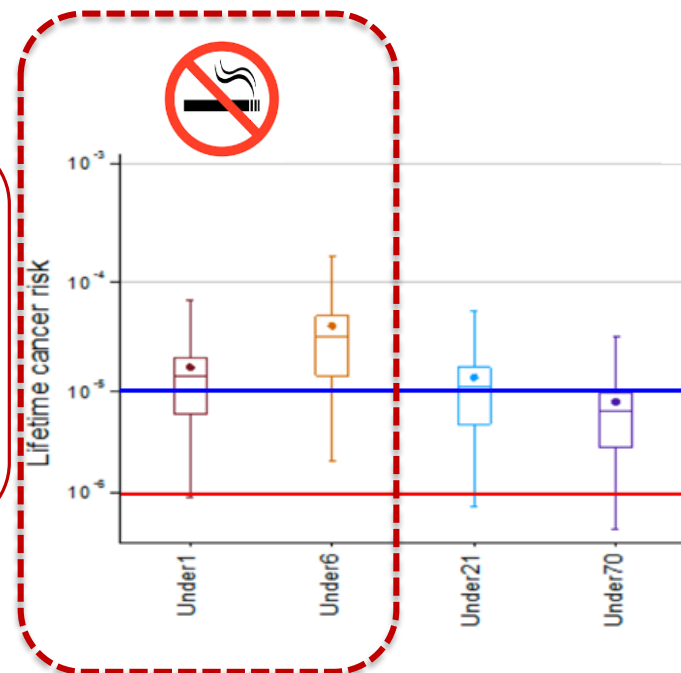


Lifetime cancer risk by age group.

Threshold values

WHO:  $10^{-5}$

USEPA:  $10^{-6}$







**SMOKER'S HOMES**  
(n=22)

Nicotine 26  $\mu\text{g/g}$   
TSNAs 0.5  $\mu\text{g/g}$



**NON-SMOKER'S  
HOMES (n=24)**

Nicotine 3.3  $\mu\text{g/g}$   
TSNAs 0.09  $\mu\text{g/g}$

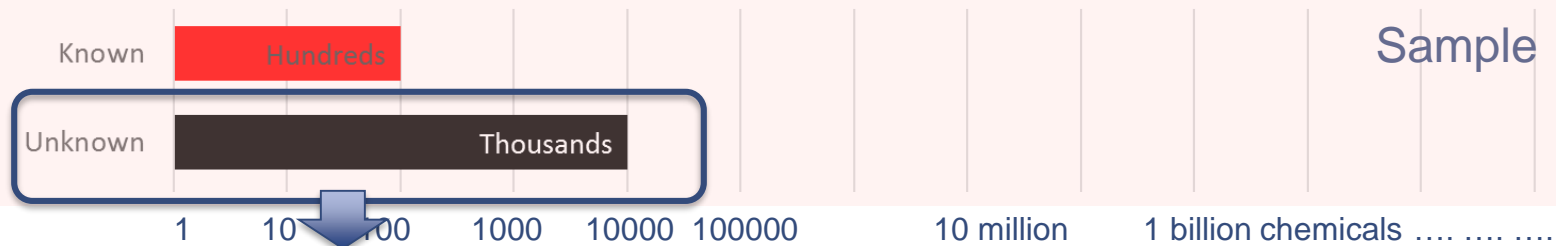
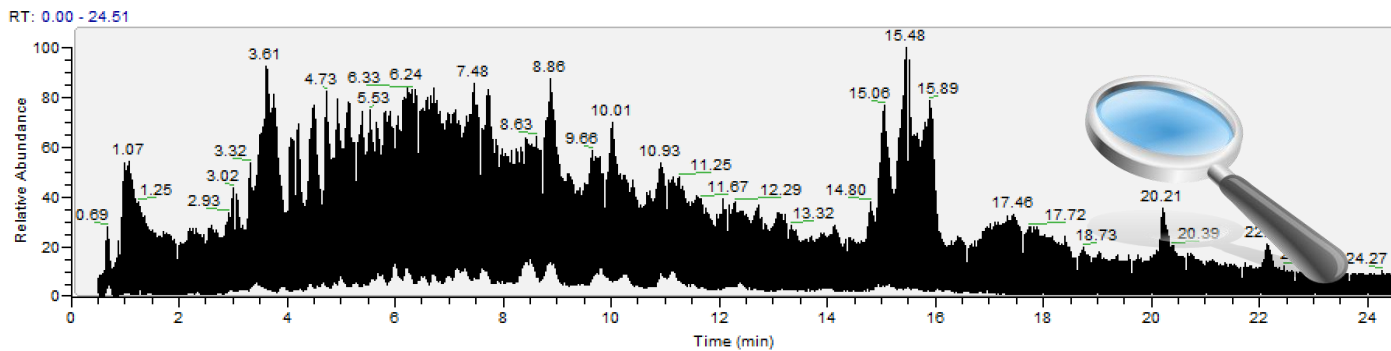
## Motivation for NTS: What else is in there?

THS-RESEARCH TEAM

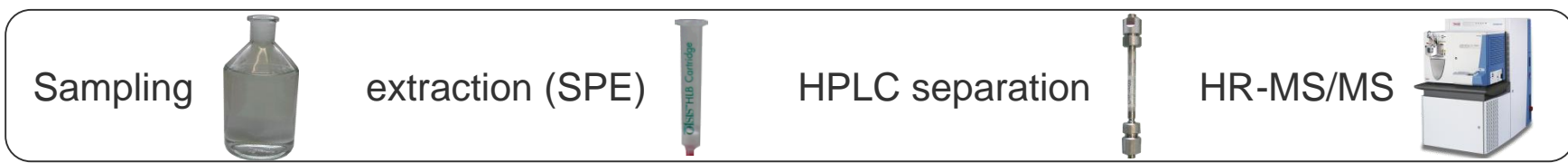


# The Challenge in Identifying Chemicals

High resolution  
mass spectrometry  
+  
Cheminformatics  
approaches  
=  
ECI@LCSB



Category	Quantity	Source/Icon
Suspected Neurotoxicants	Hundreds	
Mass Spectral Libraries	~20,000	MassBank.eu
CompTox Dashboard	~760,000	EPA Chemicals
PubChem Compound	>96 million	PubChem
1st Gen. PubChem Metabolites	>2 billion	PubChem
Generated Structures	Millions of billions	



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

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

Target List

**TARGET ANALYSIS**

(enviMass, vendor software)

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

**SUSPECT SCREENING**

Gather evidence (nontarget, ReSOLUTION, RMassBank)

**NON-TARGET SCREENING**

Componentization (nontarget)

Prioritization (enviMass)

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

# Terminology: Target, Suspect, Non-target/Unknown

---

## **TARGET:**

Known chemical, reference standard available *in house*

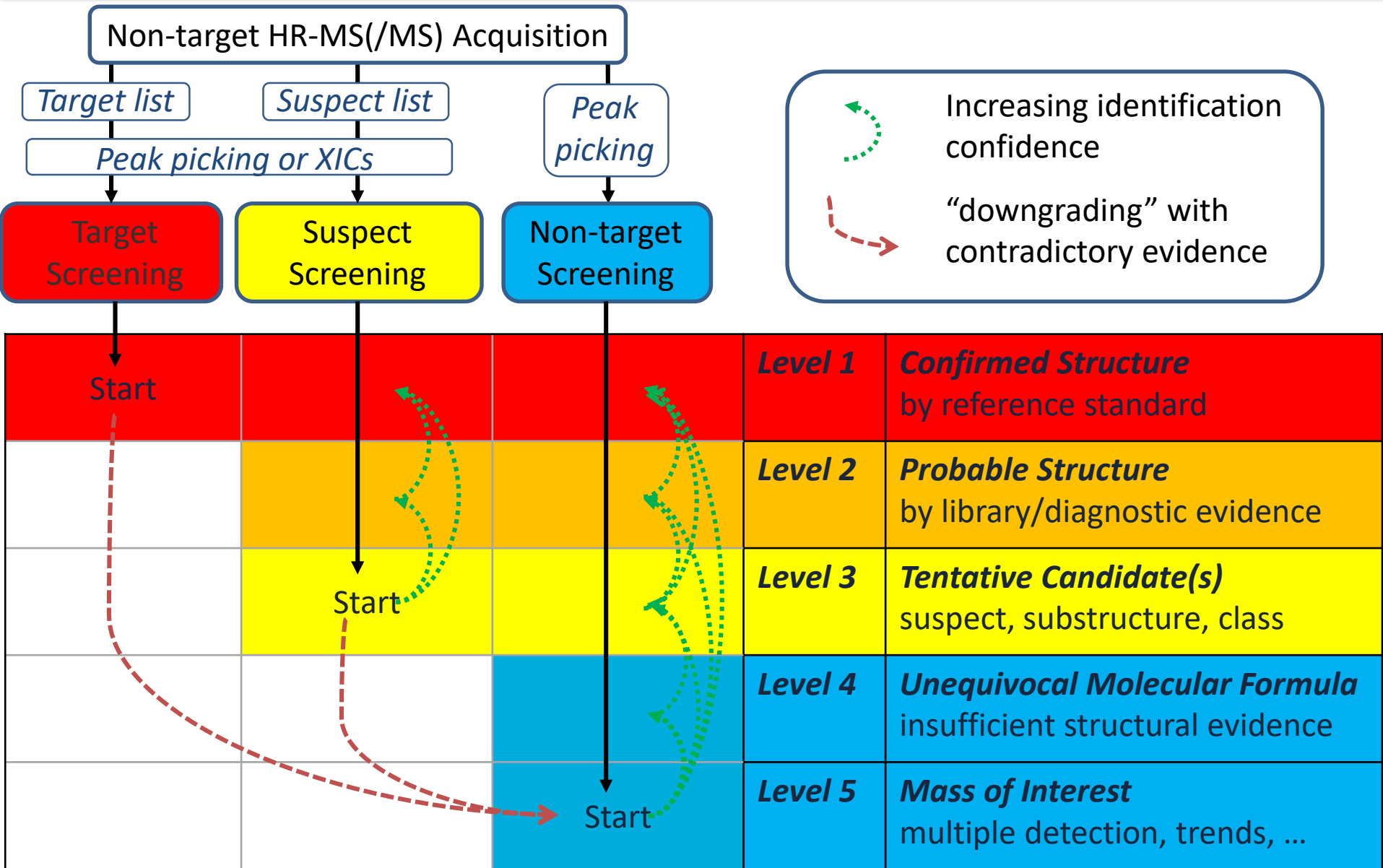
## **SUSPECT:**

Chemical suspected to be present in the sample, std. not (necessarily) available

## **NON-TARGET/ UNKNOWN:**

Mass/feature of interest detected in the sample, identity unknown

# Identification Strategies and Confidence



# THS & Target, Suspect, Non-target/Unknown



## TARGET:

Known chemical, reference standard available *in house*

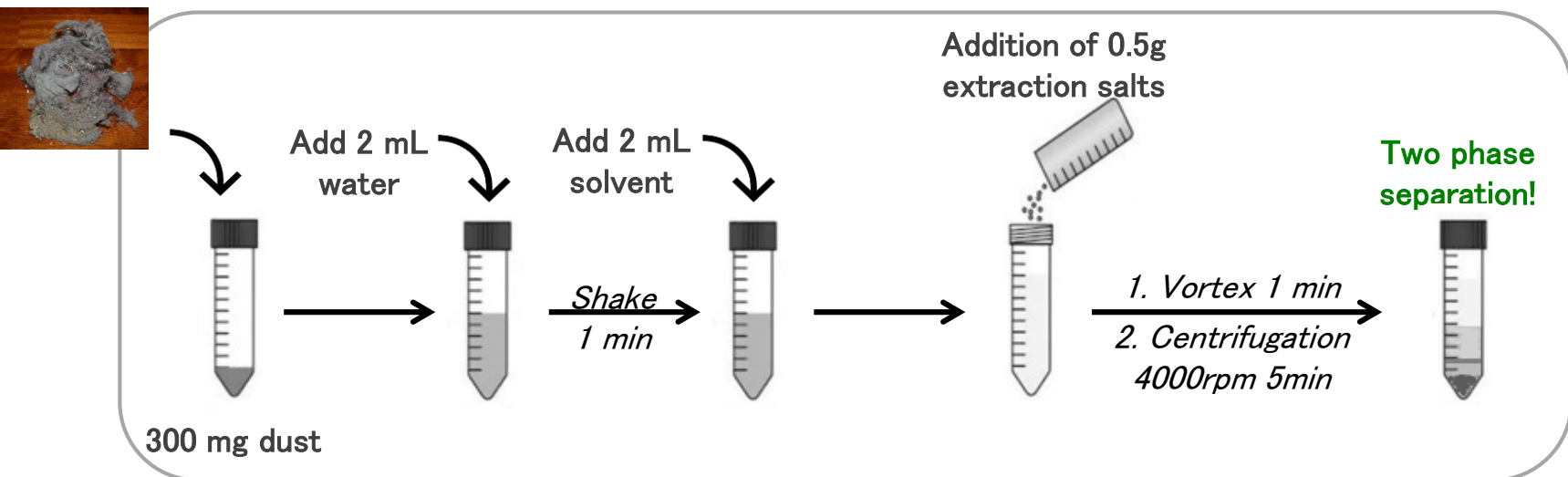
## SUSPECT:

Chemical suspected to be present in the sample, std. not (necessarily) available

## NON-TARGET/ UNKNOWN:

Mass/feature of interest detected in the sample, identity unknown

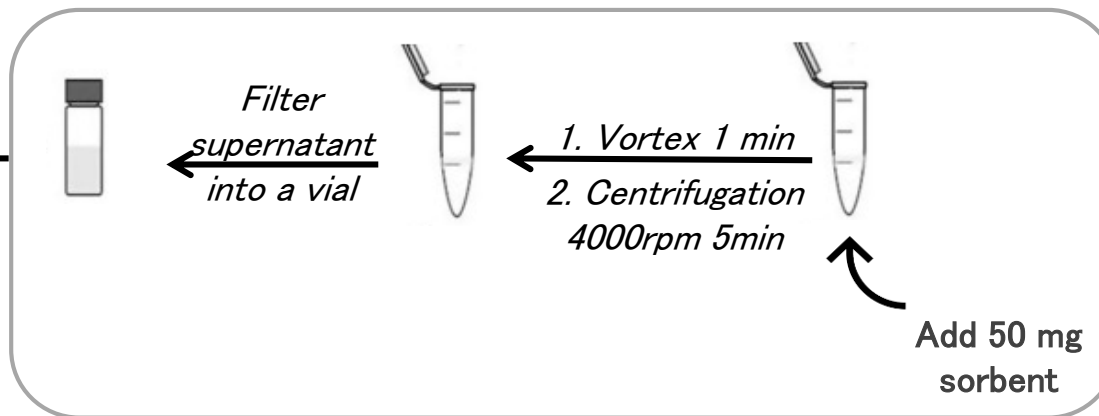
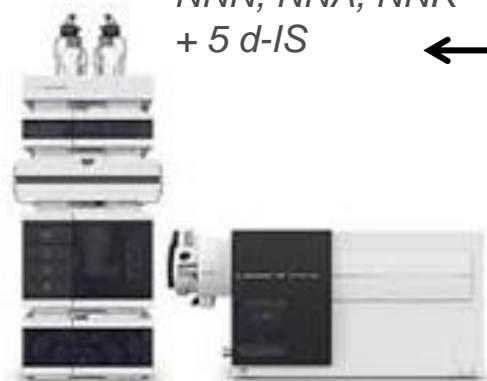
# Extraction (QuEChERS) and Target Analysis



Transfer supernatant

## UHPLC-QQQ Analysis

Nicotine, Cotinine, NNAL,  
NNN, NNA, NNK  
+ 5 d-IS



**Target results essential to verify the non-target methods!**

# THS & Target, Suspect, Non-target/Unknown



## TARGET:

Known chemical, reference standard available *in house*

## SUSPECT:

Chemical suspected to be present in the sample, std. not (necessarily) available

## NON-TARGET/ UNKNOWN:

Mass/feature of interest detected in the sample, identity unknown



# Compiled THS-specific Suspect List (n=95)



	Compound	Chemical name	CAS No.	EPA ID	Molecular Formula	Monoisotopic Mass
<b>Tobacco Specific Nitrosamines TSNA's</b>	Nicotine	3-(1-methylpyrrolidin-2-yl)pyridine	54-11-5	DTXSID1020930	C10H14N2	162.115698
	Cotinine	1-methyl-5-(pyridin-3-yl)pyrrolidin-2-one	486-56-6	DTXSID1047576	C10H12N2O	176.094963
	NNN	N'-Nitrosoanornicotine	16543-55-8	DTXSID4021476	C9H11N3O	177.090212
	3-hydroxycotinine	3-Hydroxy-1-methyl-5-(3-pyridinyl)-2-pyrrolidinone	34834-67-8	DTXSID30873224	C10H12N2O2	192.089878
	NNA	4-(methylnitrosamino)-4-(3-pyridyl)butanal	64091-90-3	DTXSID00897139	C10H13N3O2	207.233
	NNK	4-(Methylnitrosoamino)-1-(3-pyridinyl)-1-butanone	64091-91-4	DTXSID3020881	C10H13N3O2	207.100777
	NNAL	4-(methylnitrosamino)-1-(3-pyridyl)-1-butanol	76014-81-8	DTXSID8020880	C10H15N3O2	209.116427
	NAB	N'-nitrosoanabasine	37620-20-5	DTXSID3021019	C10H13N3O	191.23
	NAT	N'-nitrosoanatabine	887407-16-1	DTXSID40868005	C10H11N3O	189.21
<b>Nicotine Alkaloids</b>	Myosmine	3-(3,4-Dihydro-2H-pyrrol-5-yl)-pyridine	532-12-7	DTXSID70891866	C9H10N2	146.19
	$\beta$ -nicotyrine	Pyridine, 3-(1-methyl-1H-pyrrol-2-yl)-	487-19-4	DTXSID3075048	C10H10N2	158.084398
	2,3'-Bipyridine	2,3'-Bipyridine	581-50-0	DTXSID00206823	C10H8N2	156.068748
	N-Formylnornicotine	2-(Pyridin-3-yl)pyrrolidine-1-carbaldehyde	3000-81-5	DTXSID30336006	C10H12N2O	176.094963
	Nicotelline	Nicotelline	494-04-2	DTXSID40197781	C15H11N3	233.095297
<b>Secondary Products of Nicotine heterogeneous nitrosation</b>	Methyl nicotinate	3-Pyridinecarboxylic acid, methyl ester	93-60-7	DTXSID7044471	C7H7NO2	137.047678
	N-methylnicotinamide	N-Methylpyridine-3-carboxamide	114-33-0	DTXSID00870467	C7H8N2O	136.15

- Add these to “inclusion list” for DDA-MS/MS
- Now added to CompTox Dashboard & NORMAN-SLE & Zenodo

## NORMAN Suspect List Exchange


### Recent uploads

May 6, 2019 (NORMAN-SLE-S52.0.1.0)

Dataset

Open Access

#### S52 | THSMOKE | Thirdhand Smoke (THS) Compounds

Torres, Sonia;  Schymanski, Emma; Ramirez, Noelia;

This is the collection associated with list S52 THSMOKE on the NORMAN Suspect List Exchange. <https://www.norman-network.com/?q=suspect-list-exchange> S52 THSMOKE T (THS) Compounds THSMOKE XLSX , CSV (06/05/2019) CompTox THSMOKE I InChIKeys (06/05/

Uploaded on May 6, 2019

March 1, 2018 (NORMAN-SLE-S0.0.1.0)

Dataset

Open Access

#### S0 | SUSDAT | Merged NORMAN Suspect List: SusDat

91

 views

55

 downloads

[See more details...](#)



Tweeted by 5

[See more details](#)

[New upload](#)

### Suspect List

Public repository (under license) for suspect lists currently hosted in the NORMAN Suspect List Exchange. <https://www.norman-network.com/?q=suspect-list-exchange>

[Read more](#)

**Publication date:**

May 6, 2019

**DOI:**

DOI [10.5281/zenodo.2669467](https://doi.org/10.5281/zenodo.2669467)

**Keyword(s):**

[Suspect Screening](#) [Thirdhand smoke](#)

**Communities:**

[NORMAN Suspect List Exchange](#)

**License (for files):**

[Creative Commons Attribution 4.0 International](#)

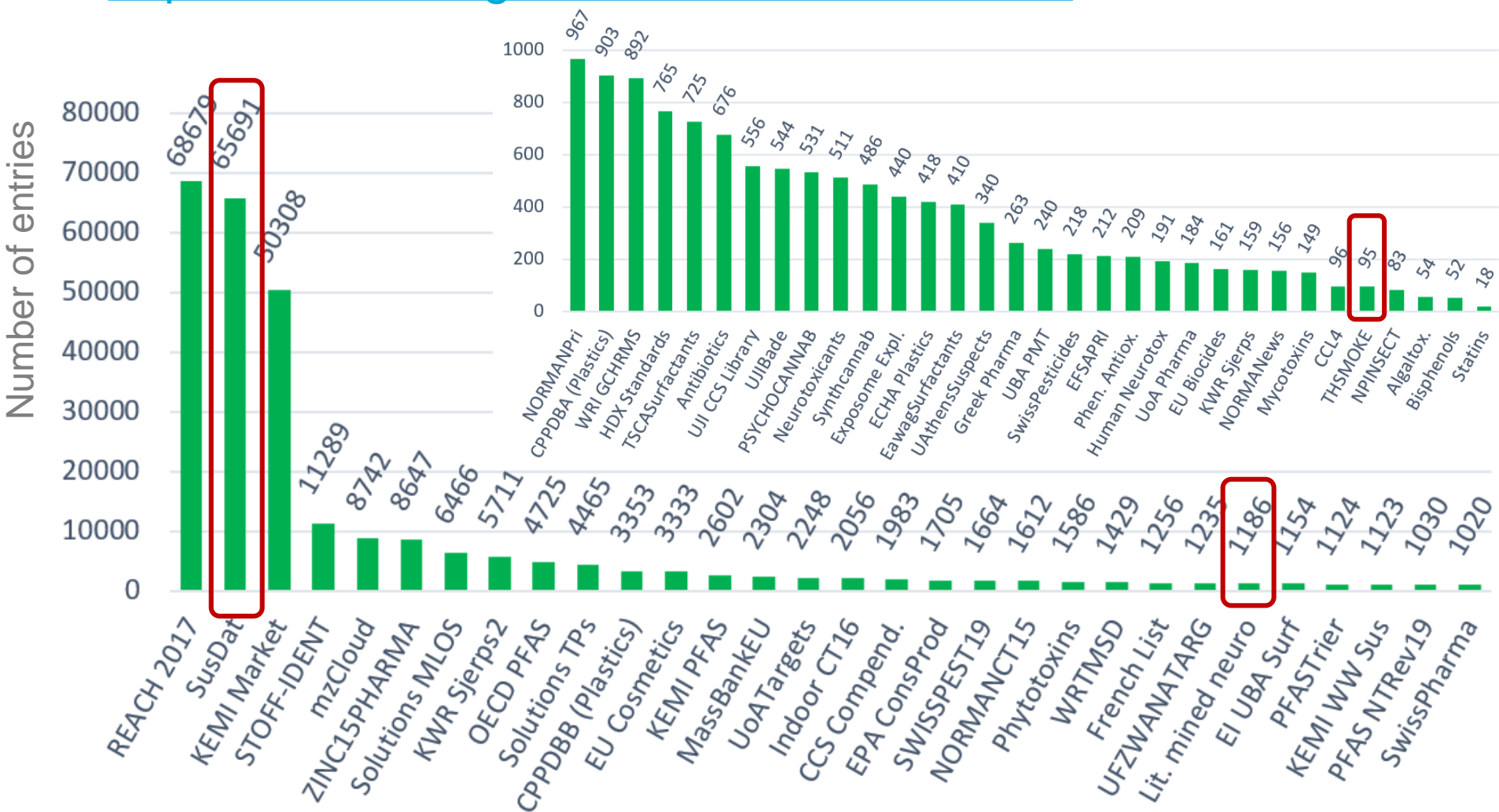
**Policy:**

This community will collect data that is relevant and public (or to be public) for

# NORMAN Suspect List Exchange (SLE)



- <https://www.norman-network.com/nds/SLE/> ...now 62 lists!
- <https://zenodo.org/communities/norman-sle> ... with DOI



# THS & Target, Suspect, Non-target/Unknown



## TARGET:

Known chemical, reference standard available *in house*

## SUSPECT:

Chemical suspected to be present in the sample, std. not (necessarily) available

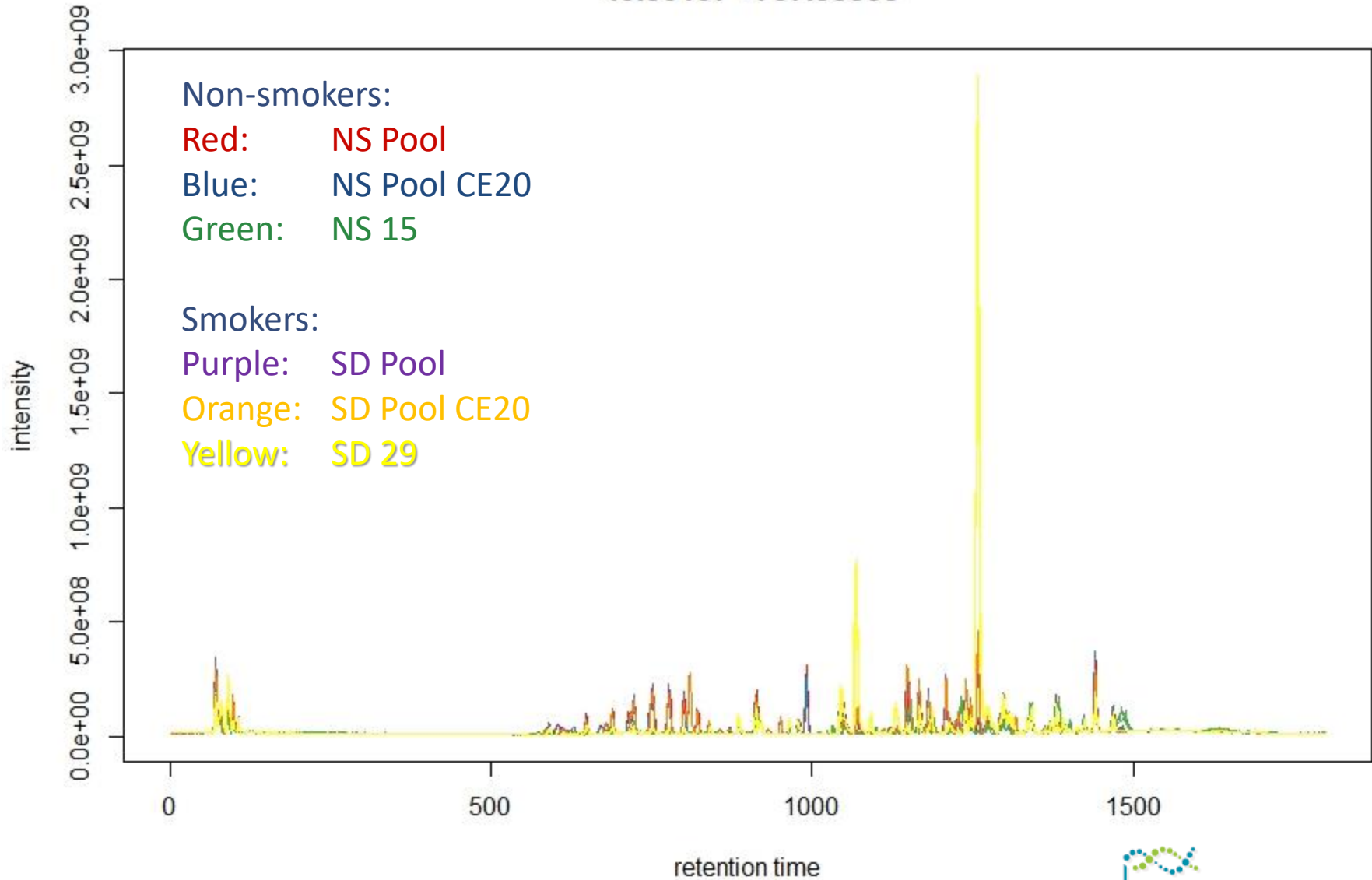
## NON-TARGET/ UNKNOWN:

Mass/feature of interest detected in the sample, identity unknown

# IISPV-URV XCMS-based Metabolomics Workflow



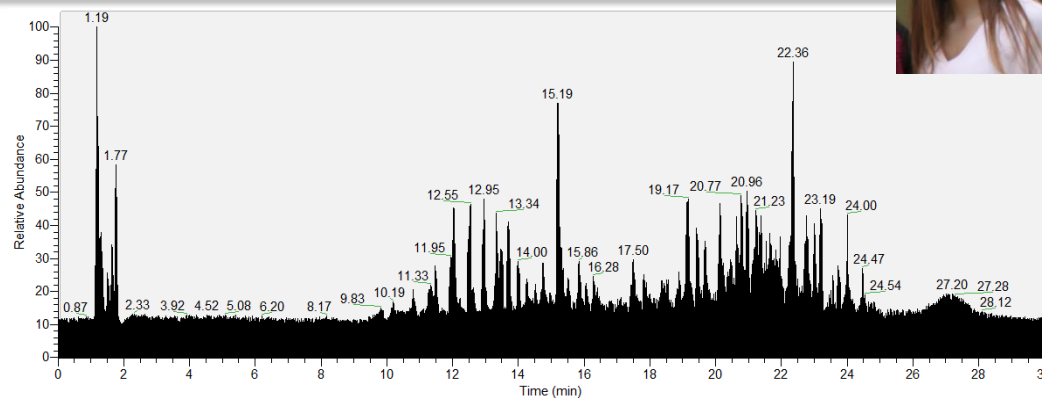
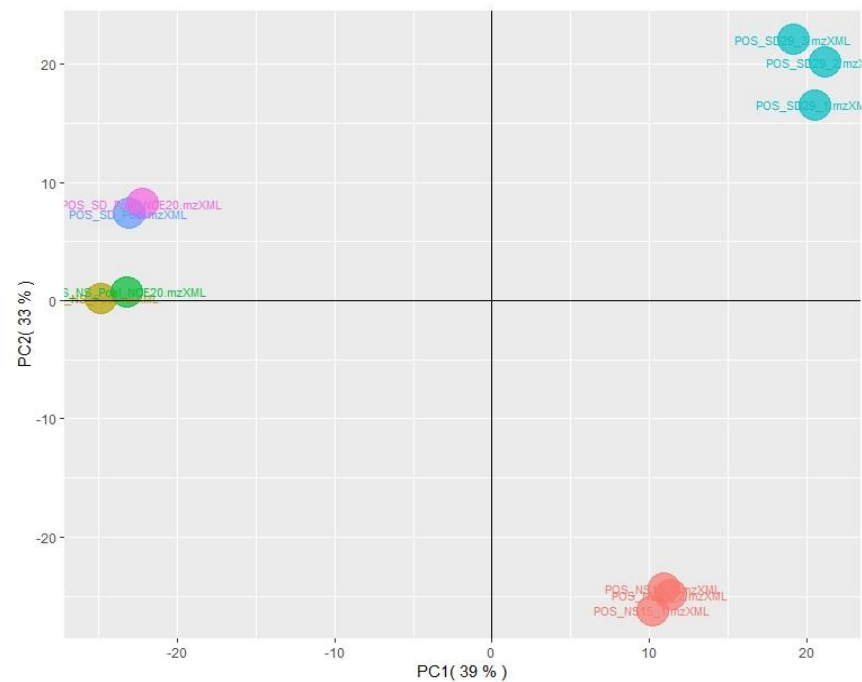
49.50197 - 757.55888



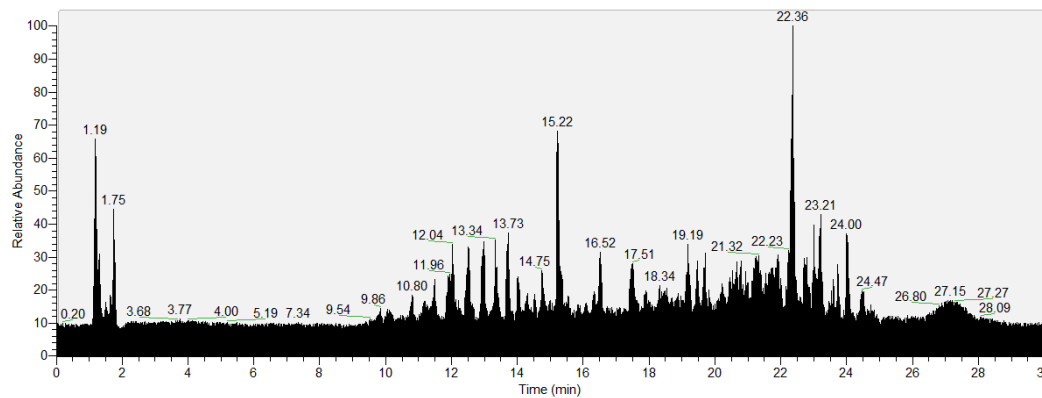
# IISPV-URV XCMS-based Metabolomics Workflow



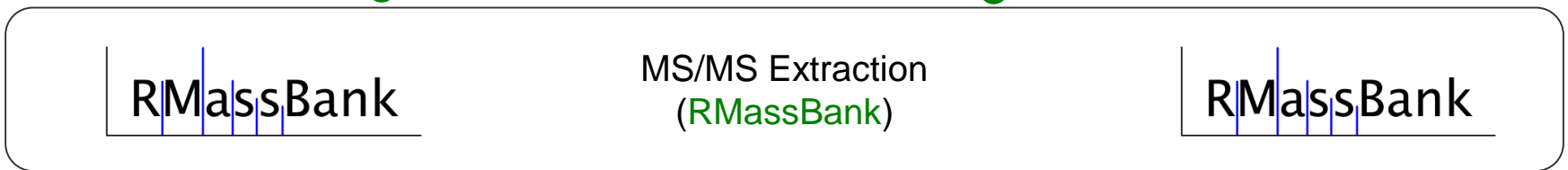
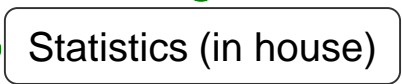
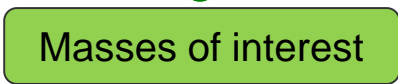
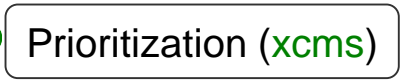
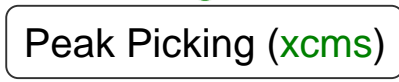
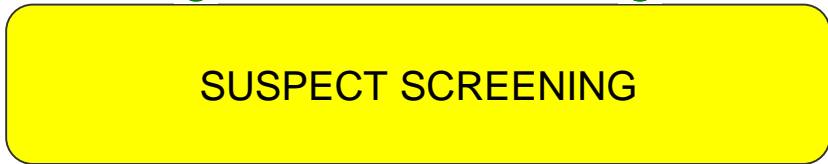
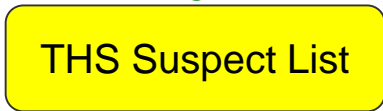
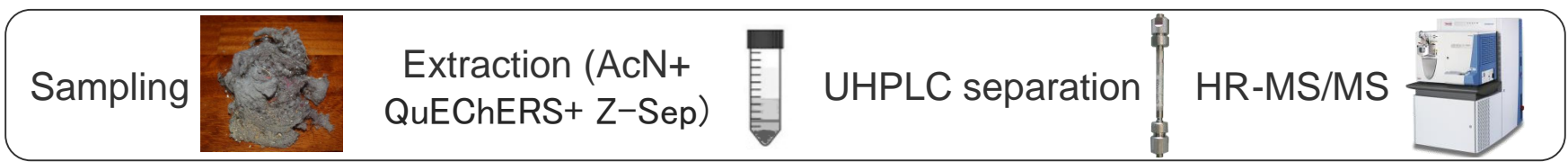
class

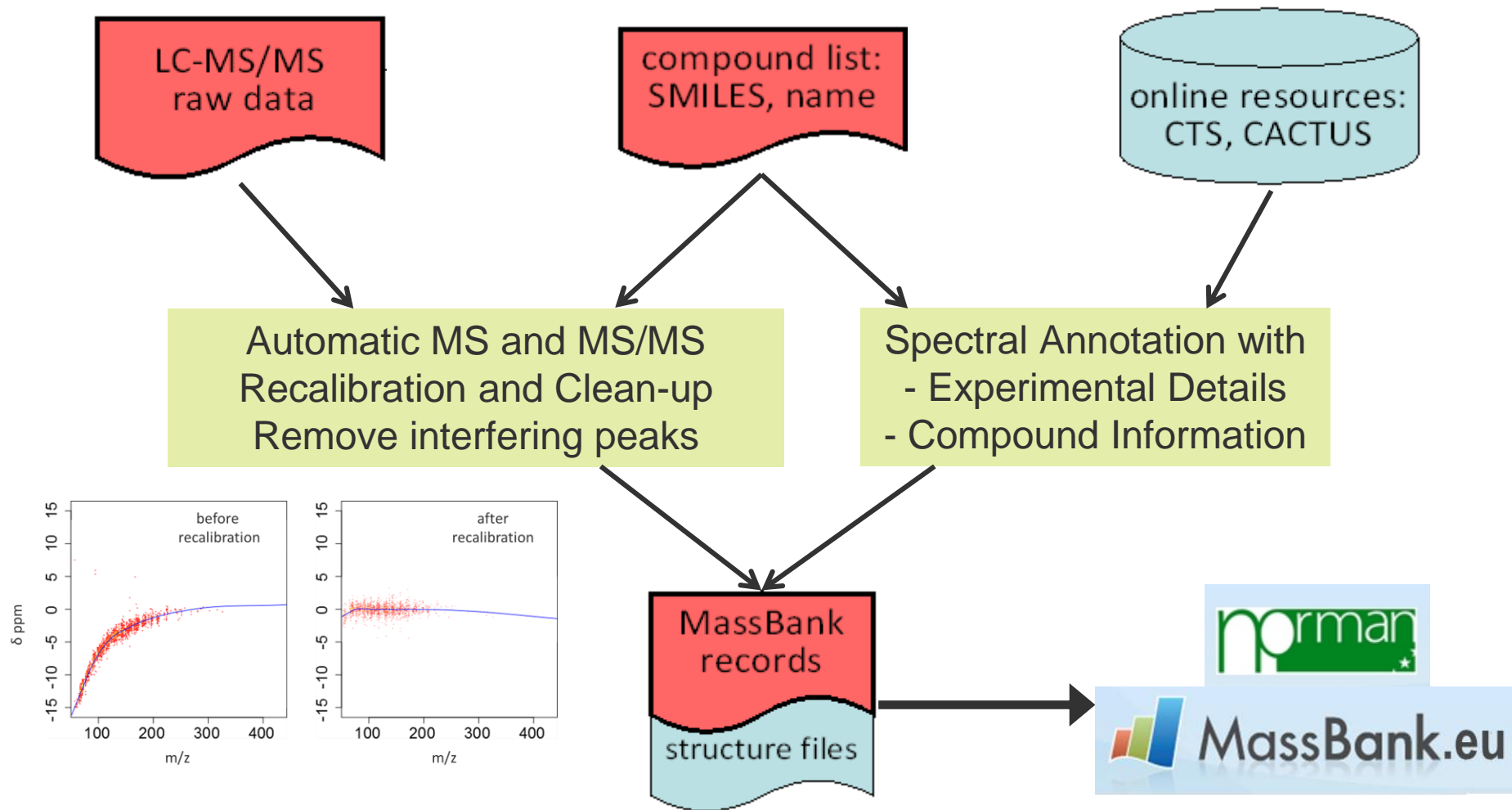


Smoker Pool



Non Smoker Pool

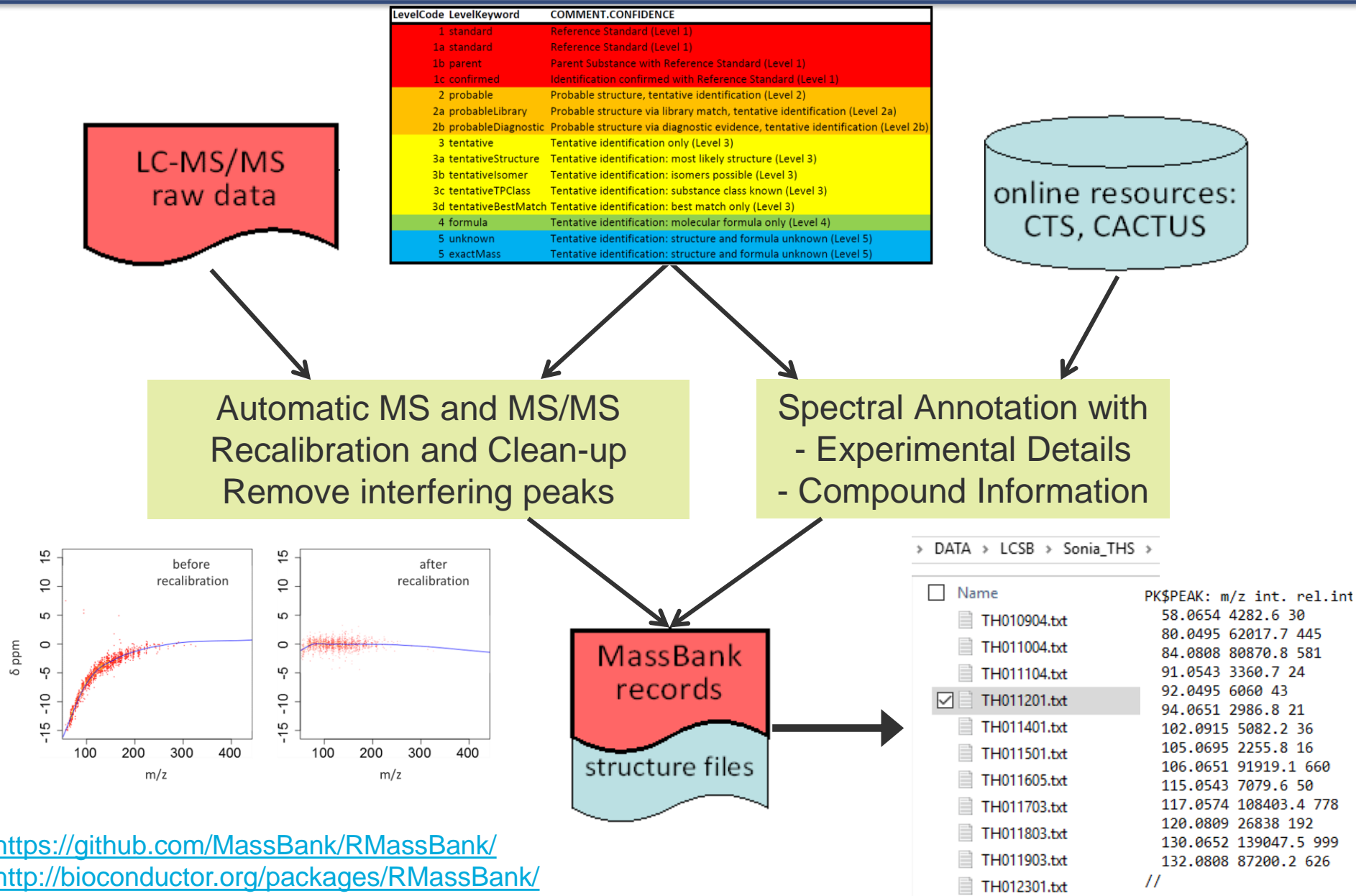




<https://github.com/MassBank/RMassBank/>  
<http://bioconductor.org/packages/RMassBank/>



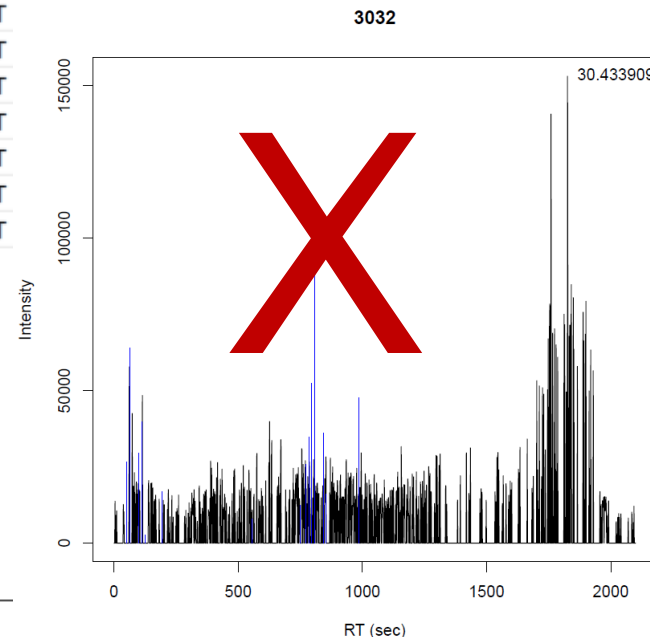
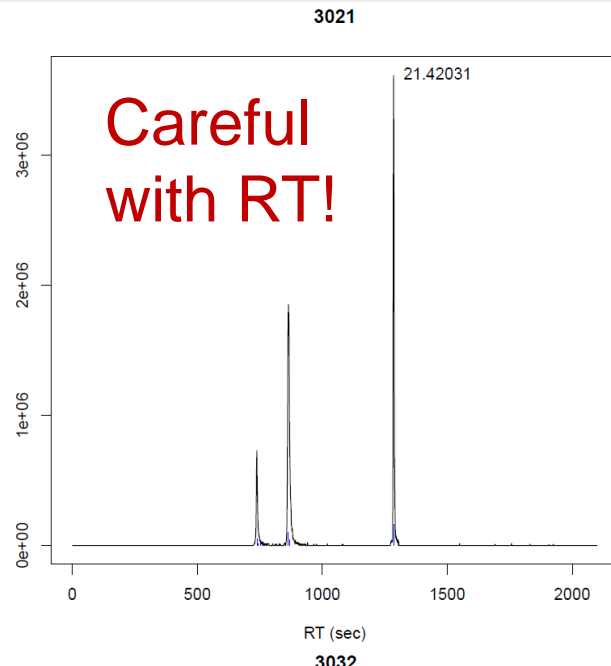
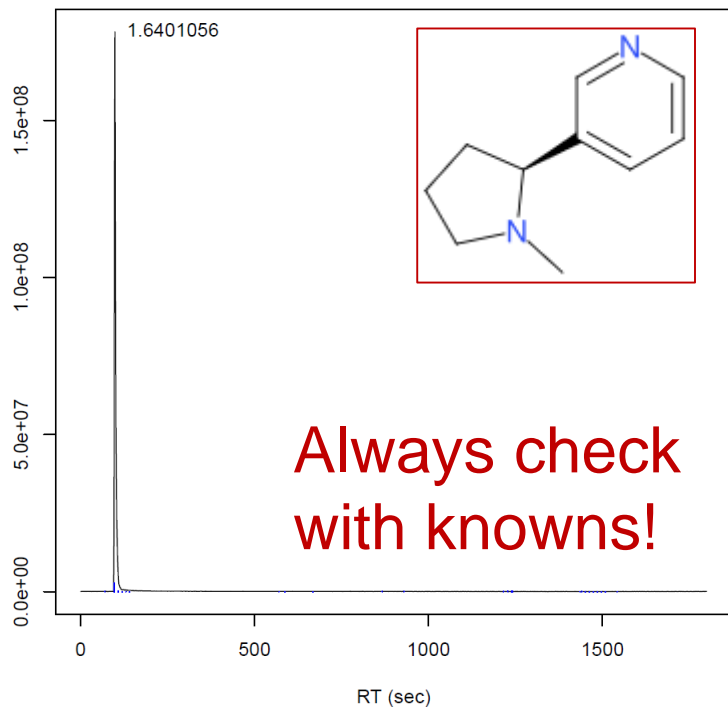
# Extracting Mass Spectra for NTS



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

# MS/MS Extraction with RMassBank

A	B	C	D	E	F	G	H	I	
ID	Name	Name2	DTXSID	SMILES	CAS	RT	Formula	mz	Std
112	Nicotine (3	Nicotine	DTXSID102	CN1CCC[C@H]1N	54-11-5	1.62	C10H14N2		Std
114	Nicotine-d3	DL-Nicotine	DTXSID804	[2H]C([2H])[2H]1CCN1	69980-24-1	1.62	C10H11D3N2		Std
117	Cotinine (1	Cotinine	DTXSID104	CN1[C@@H](C(=O)N1C)C	486-56-6	1.77	C10H12N2O		Std
119	NNN (N'-N'-Nitrosor		DTXSID402	O=NN1CCC1	16543-55-8	1.79	C9H11N3O		Std
122	Cotinine-d3	Cotinine-d3	DTXSID404	[2H]C([2H])[2H]1CCN1C(=O)N1	110952-70	1.77	C10H9D3N2O		Std
134	NNK (4-(M	4-(N-Methyl	DTXSID302	CN(CCCC(=O)N1C)C1	64091-91-4	11.33	C10H13N3O2		Std
135	NNA (4-(m	4-(methyl	DTXSID008	CN(N=O)C1CCC1	64091-90-3	1.77	C10H13N3O2		Std
136	NNAL (4-(n	4-(Methyl	DTXSID802	CN(CCCC(O)N1C)C1	76014-81-8	1.77	C10H15N3O2		Std
157	isoNNAL			CN(C(CCCO)N1C)C1		1.77	C10H15N3O2		Std
158	NNN				112		7D4N3O		Std
159	NNA						H10D3N3O2		Std
160	NNA						H12D3N3O2		Std
1001	FTO						60.04523	NT	
1002	FTO						61.04228	NT	
1003	FTO						90.97725	NT	
1004	FTO						104.1072	NT	
1005	FTO						109.1013	NT	
1006	FTO						111.0216	NT	
1007	FTO						114.0666	NT	
1008	FTO						121.0398	NT	
1009	FTO						125.0392	NT	
1010	FTO						125.0959	NT	



# Coming soon ... pre-screening with ShinyScreen!

Activities GNU Icecat Web Browser 17 Jan 11:17 AM ShinyScreen - GNU IceCat

ShinyScreen

127.0.0.1:5254

Inputs

- Compound list
- Compound sets

Plot

741 EIC (mz= 296.116031008)

Retention time at max. intensity (MS1)

- 15 ; rt= 20.44 min
- 30 ; rt= 20.44 min
- 45 ; rt= 20.44 min
- 60 ; rt= 20.44 min
- 75 ; rt= 20.43 min
- 90 ; rt= 20.43 min

CC(C)(O)C(Oc1ccc(Cl)cc1)n2cncn2

MS2

Retention time at max. intensity (MS2)

- 15 ; rt= 20.44 min

MS2

Retention time at max. intensity (MS2)

- 15 ; rt= 20.44 min

Prescreening analysis

Quality Control

- MS1
- MS2
- Alignment
- AboveNoise

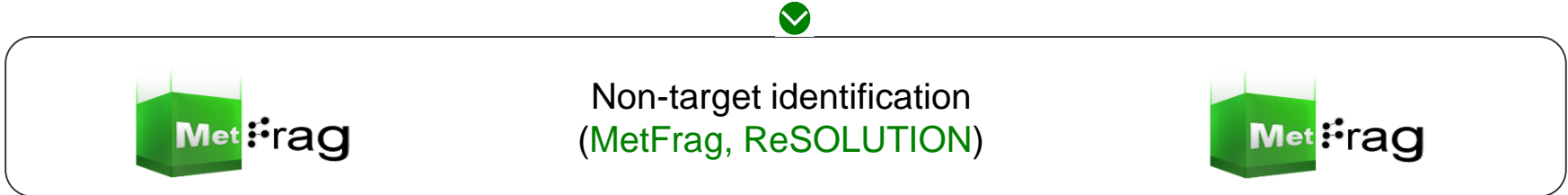
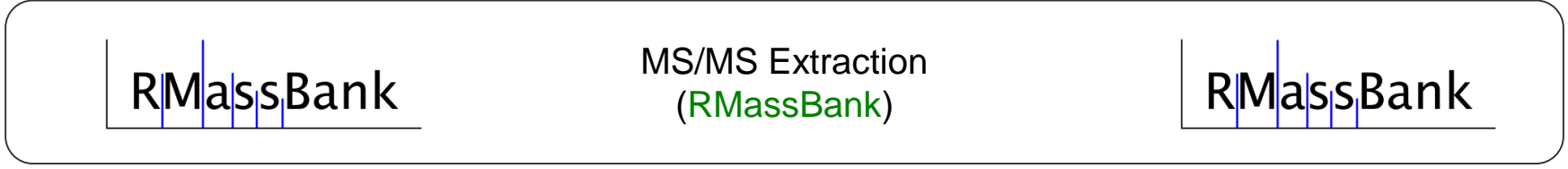
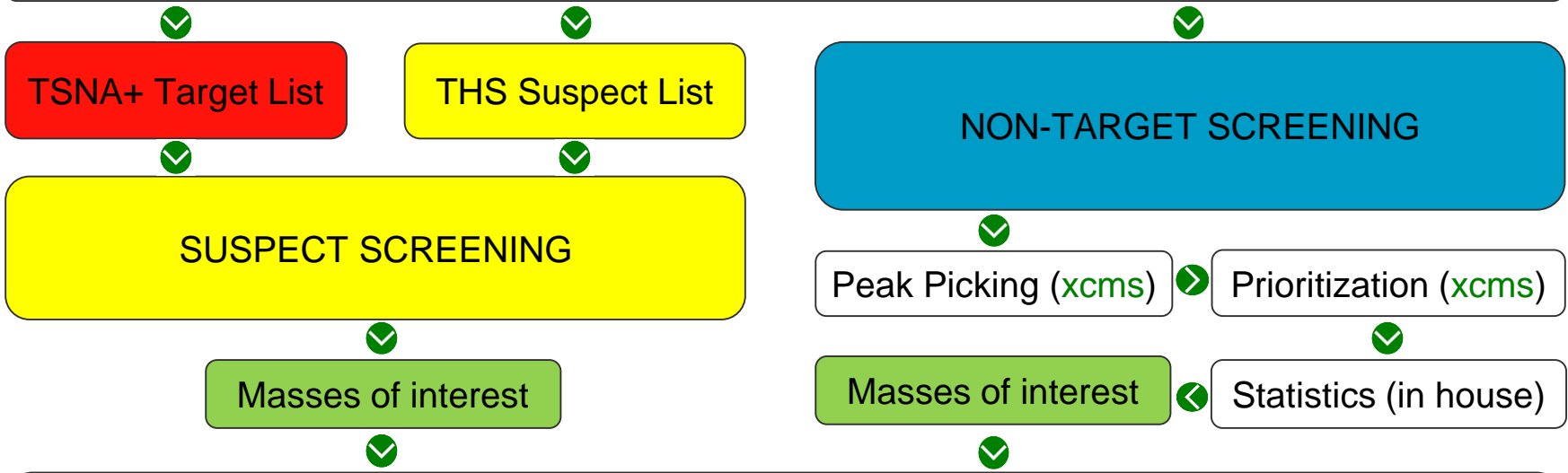
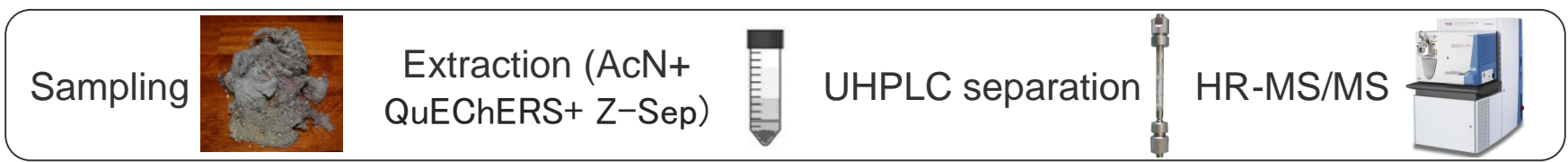
Comments:

Insert your comment here...

Todor Kondic, Mira Narayanan,  
Jessy Krier, Anjana Elapavalore,  
Hiba Mohammed Taha.

Figure source: Todor Kondic





# Tentative Identification with MetFrag

Web Interface: <http://msbi.ipb-halle.de/MetFrag/>



MetFrag

In silico fragmentation for computer assisted identification of metabolite mass spectra

## Database Settings

Database:

Neutral Mass:  Search ppm:

Formula:

Identifiers:

Parent Ion:

## Candidate Filter & Score Settings

### Fragmentation Settings & Processing

Mzppm:

Mzabs:

Mode:

#### MS/MS Peak list

90.97445 681  
106.94476 274  
110.02750 110  
115.98965 95  
117.98540 384  
124.93547 613  
124.99015 146  
125.99793 207  
133.95592 777  
143.98846 478  
144.99625 352  
.....

# Tentative Identification with MetFrag

Web Interface: <http://msbi.ipb-halle.de/MetFrag/>



## MetFrag

In silico fragmentation for computer assisted identification of metabolite mass spectra

### Database Settings

Database:

Neutral Mass:  Search

Formula:

Identifiers:

### Candidate Filter & Score Settings

### Candidate Filter & Score Settings

#### Candidate Filters:

- Element Inclusion
- Element Exclusion
- Substructure Inclusion
- Substructure Exclusion
- Substructure Information
- Minimum Number Elements

#### MetFrag Scoring Terms:

- Substructure Inclusion
- Substructure Exclusion
- Retention Time
- Suspect Inclusion Lists
- Spectral Similarity (MoNA)
- Exact Spectral Similarity (MoNA)

Mzppm:

Mzabs:

Mode:

Parent Ion:



MetFrag

In silico fragmentation for computer assisted identification of metabolite mass spectra

## Database Settings

Database:  Include references:

Neutral Mass:  Search ppm:

Formula:

Identifiers:

Parent Ion:

<https://comptox.epa.gov/dashboard/>



MetFrag

In silico fragmentation for computer assisted identification of metabolite mass spectra

## Database Settings

Database:

CompTox\_01May18\_Select

Neutral Mass:

Formula:

Identifiers:

CSV

PSV

SDF

### Local Databases

CompTox\_01May18\_AllMetaData

CompTox\_01May18\_SelectMetaData

CompTox\_01May18\_SelectMetaDataPlu

Parent Ion:

163.123

[M+H]<sup>+</sup>

Calculate

Retrieve Candidates

Download Candidates

Candidate Filter & Score Settings

Retrieve Candidates



187 Candidates



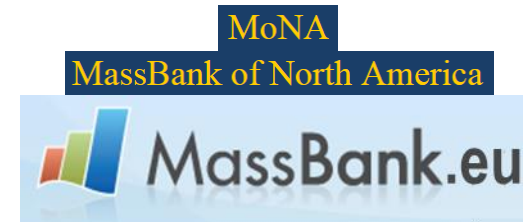
## Candidate Filter & Score Settings

### Candidate Filters

- Element Inclusion
- Element Exclusion
- Substructure Inclusion
- Substructure Exclusion
- Substructure Information
- Minimum Number Elements

### MetFrag Scoring Terms

- Substructure Inclusion
- Substructure Exclusion
- Retention Time
- Suspect Inclusion Lists
- Spectral Similarity (MoNA)
- Exact Spectral Similarity (MoNA)
- Statistical Scoring



MASSBANKEU  
 NORMANSUSDAT  
 NUMBER\_OF\_PUBMED\_ARTICLES  
 PUBCHEM\_DATA\_SOURCES  
 TOX21SL  
 TOXCAST  
 TOXCAST\_PERCENT\_ACTIVE

Select Item(s) 4 of 11 item(s) selected

### Database Scoring Terms

Select Item(s) 4 of 11 item(s) selected

## ○ Fragmentation Settings and Processing

### Fragmentation Settings & Processing

Mzppm:

Mzabs:

Mode:

Tree depth:

Group candidates

Process Candidates

MS/MS Peak list

58.0654	4282.6	30
80.0495	62017.7	445
84.0808	80870.8	581
91.0543	3360.7	24
92.0495	6060	43
94.0651	2986.8	21
102.0915	5082.2	36
105.0695	2255.8	16
106.0651	91919.1	660
115.0543	7079.6	50
117.0574	108403.4	

Show Spectrum

Download Parameters

```
PK$PEAK: m/z int. rel.int
58.0654 4282.6 30
80.0495 62017.7 445
84.0808 80870.8 581
91.0543 3360.7 24
92.0495 6060 43
94.0651 2986.8 21
102.0915 5082.2 36
105.0695 2255.8 16
106.0651 91919.1 660
115.0543 7079.6 50
117.0574 108403.4 778
120.0809 26838 192
130.0652 139047.5 999
132.0808 87200.2 626
//
```

# MetFrag – Example: CompTox + Nicotine V

## Results

**Weights**

MetFrag (1st)	<input type="range" value="100"/>	100 %
ExactSpectralSimilarity (2nd)	<input type="range" value="100"/>	100 %
DATA_SOURCES (3rd)	<input type="range" value="100"/>	100 %
NORMANSUSDAT (4th)	<input type="range" value="100"/>	100 %
NUMBER_OF_PUBMED_ARTICLES (5th)	<input type="range" value="100"/>	100 %
TOXCAST_PERCENT_ACTIVE (6th)	<input type="range" value="100"/>	100 %

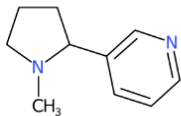
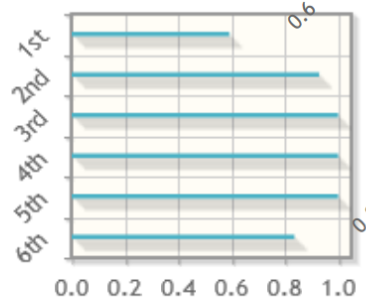
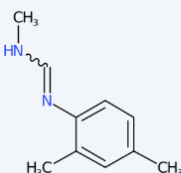
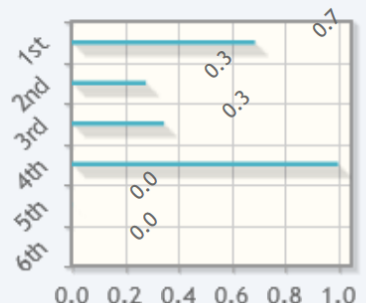
Download Results

Filter Candidates by explained MS/MS Peaks

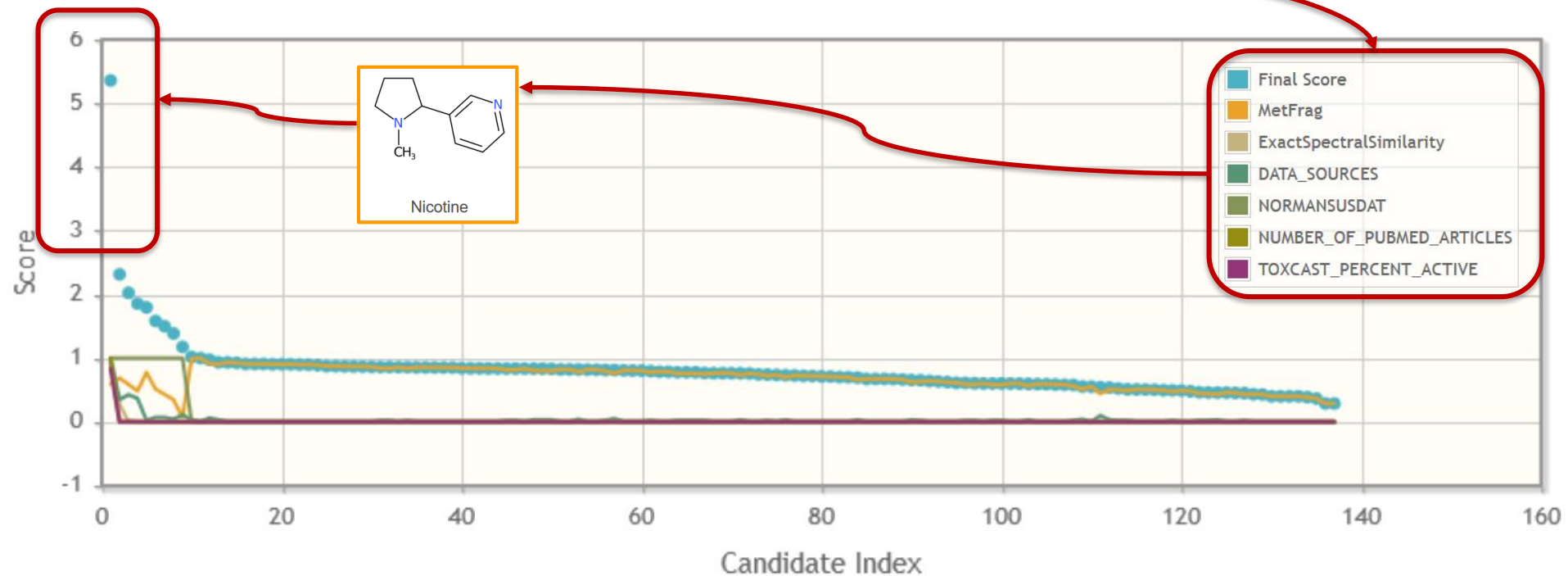
MS/MS Peaks

Filter Candidates

Navigation: << 1 2 3 4 5 6 7 8 9 10 >>

#	Molecule	Identifier	Mass	Formula	Normalized Scores	FinalScore	Details
1	 <p>Nicotine</p>	<p><a href="#">DTXSID1020930</a></p> <p> <a href="#">DTXSID8021725</a>  <a href="#">DTXSID3048154</a>  <a href="#">DTXSID0046351</a>  <a href="#">DTXSID6020931</a>  <a href="#">DTXSID5075319</a>  <a href="#">DTXSID3088421</a> </p> <p>InChIKeyBlock1 = <a href="#">SNICXCGAKADSCV</a></p>	162.11576	C <sub>10</sub> H <sub>14</sub> N <sub>2</sub>		5.3553	<p>Peaks: 10 / 14</p> <p>Fragments</p> <p>Scores</p> <p>Download</p>
2	 <p>N'-(2,4-Dimethylphenyl)-N-methylformamidi</p>	<p><a href="#">DTXSID1037696</a></p> <p><a href="#">DTXSID10199510</a></p> <p>InChIKeyBlock1 = <a href="#">JIIOLEGNERQDIP</a></p>	162.11576	C <sub>10</sub> H <sub>14</sub> N <sub>2</sub>		2.3137	<p>Peaks: 9 / 14</p> <p>Fragments</p> <p>Scores</p> <p>Download</p>

# Connecting Resources in MetFrag



- Create config files (MetFragConfig)
- Run MetFrag (runMetFrag)

```
176 # now, run MetFrag and extract results for reporting into compd_info.
177
178 results_filename <- paste0(run_name,"_",compdID_char,"_",as.character(i))
179 if (isPos) {
180     config_file <- MetFragConfig(mass = ExactMass, adduct_type = 0, neutralPrecursorMass=TRUE,
181                                 results_filename = results_filename,
182                                 peaklist_path = MetFrag_msms, base_dir = results_run_dir,
183                                 DB = "LocalCSV", localDB_path=localDB,useMonaIndiv = T,useMoNAMetFusion = T,
184                                 IsPosMode = TRUE,filter_by_InChIKey = F,rt_file_path=MetFrag_rt_file,rt_exp=RT)
185 } else {
186     config_file <- MetFragConfig(mass = ExactMass, adduct_type = 0, neutralPrecursorMass=TRUE,
187                                 results_filename = results_filename,
188                                 peaklist_path = MetFrag_msms, base_dir = results_run_dir,
189                                 DB = "LocalCSV", localDB_path=localDB,useMonaIndiv = T,useMoNAMetFusion = T,
190                                 IsPosMode = FALSE,filter_by_InChIKey = F,rt_file_path=MetFrag_rt_file,rt_exp=RT)
191 }
192
193 runMetFrag(config_file, MetFrag_dir, CL_name = "MetFrag2.4.4-msready-CL.jar")
194
195 results_file <- paste0(results_run_dir,"/results/",results_filename,".xls")
```

## ○ Extract results and summarize

```
198 #extract results we need:
199 MetFrag_res <- read_excel(results_file)
200
201 compd_info$num_poss_IDs[i] <- length(MetFrag_res$Score)
202 compd_info$poss_IDs[i] <- paste(MetFrag_res$Name,collapse=";")
203 compd_info$poss_ID_scores[i] <- paste(MetFrag_res$Score,collapse=";")
204 compd_info$max_Score[i] <- max(MetFrag_res$Score)
205 compd_info$n_Score_GE3p5[i] <- length(which(MetFrag_res$Score>=3.5))
206 compd_info$n_Score_GE3[i] <- length(which(MetFrag_res$Score>=3))
207 compd_info$n_Score_GE2p5[i] <- length(which(MetFrag_res$Score>=2.5))
208 compd_info$poss_DTXSIDs[i] <- paste(MetFrag_res$DTXSID,collapse=";")
209 compd_info$poss_CAS[i] <- paste(MetFrag_res$CAS,collapse=";")
210 compd_info$MoNAScore[i] <- paste(MetFrag_res$OfflineIndividualMoNAScore,collapse=";")
211 compd_info$MaxMoNAScore[i] <- max(MetFrag_res$OfflineIndividualMoNAScore)
212
213 }
214
215
216 write.csv(compd_info,paste0(results_summary_dir,"/MetFragResultSummary_",run_name,".csv"),row.names = F)
217
```

# CSV Summary Output for Results Interrogation

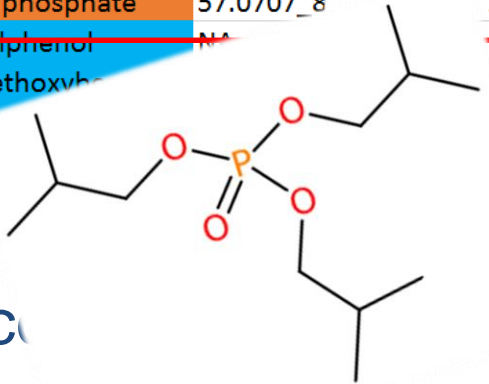
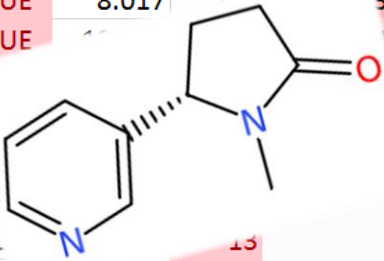
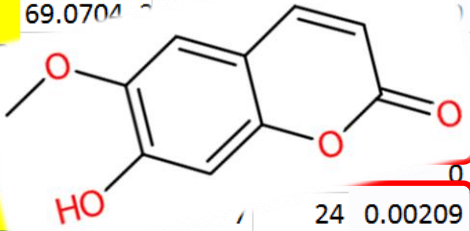
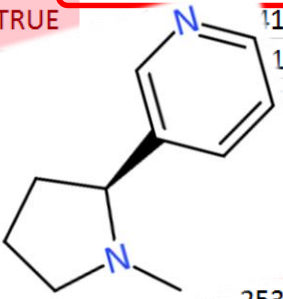
ID	mz	Name	RT	Int	MS/MS	RT_RMB	#Cand	MaxScore	SMILES	Name_maxScore	ExplPeaks	#Peaks	#Peaks	MoNA
3027	163.1221	FT0532	1.073	9E+07	TRUE	1.052	137	7.270015	CN1CCC	Nicotine	58.0658_9	9	28	0.99701
3008				+07	TRUE		11	7.22864	CN1CC(=	Creatinine	57.0454_1	3	14	0.54699
3131				+06			1	7.028138	CC(C=CC	all-trans-Retinoic acid	57.0706_5	8		7
3321				+				6.093	CCCCCC	Didecyl phthalate	69.0704			
3484				+C			6		CC(C)(C)	2,2'-Oxamidodiethyl bi...				
3206	346.1096	FT3193	21.08	3E+00				5.973483	CCCN(CC	Nitralin				
3044				+07				5.960103	CCN(CC)	N,N-Diethylnicotinamid				
3006				+06				5.94264	OCCOCC	Diethylene glycol				
3043				+07			253	5.721703	CCN(CC)	N,N-Diethylnicotinamide			24	0.00209
3055				+07	TRUE	11.202	114	5.712626	COC1=C	Scopoletin	53.0393_1	9	28	0.94271
3046	183.0796	FT0741	24.77	1E+07	TRUE	24.744	72	5.689798	O=C(C1=	Benzophenone	50.0159_2	5	45	0
3039				+07	TRUE	7.635	219	5.467559	CN1C(CC	Cotinine	53.0393_1	9	23	0.9987
3038				+07	TRUE	8.017	21		CC1CN(N	4-Methyl-1-phenylpyrazoli	53.0393_2	9	40	0.66408
3183				+07	TRUE				CCCCCC	Dihexyl phthalate	54.0452_2	9	51	0
3095				+08					CC(O)CO	Triisobutyl phosphate	57.0707_8		15	0.99613
3020	151.1111	FT0427	13.16	2E+06					CC(C)(C)	4-tert-Butylphenol				0
3029				+05	FAL				COC1=CC	1,2,4-Trimethoxyph				0
3123				+06	FAL				CCCCCC	Octadec				0
3128				+06	FAL		13		CCCCCC	Octadeca				0
3172				+06	FAL		15		CCCCCC	Stearylbe				0

Level 1  
Target

Level 2  
MSMS Match

Level 3  
Tentative ID

Level 5  
No MS/MS



⇒ Approaching automatic assignment of candidate IDs  
 ⇒ Quick, high throughput prioritization for data acquisition

Sampling



Extraction (AcN+  
QuEChERS+ Z-Sep)



UHPLC separation



HR-MS/MS



Conversion (Proteowizard)



TSNA+ Target List



THS Suspect List



NON-TARGET SCREENING



SUSPECT SCREENING



Peak Picking (xcms)



Prioritization (xcms)



Masses of interest



Masses of interest



Statistics (in house)



RMassBank

MS/MS Extraction  
(RMassBank)

RMassBank



Non-target identification  
(MetFrag, ReSOLUTION)



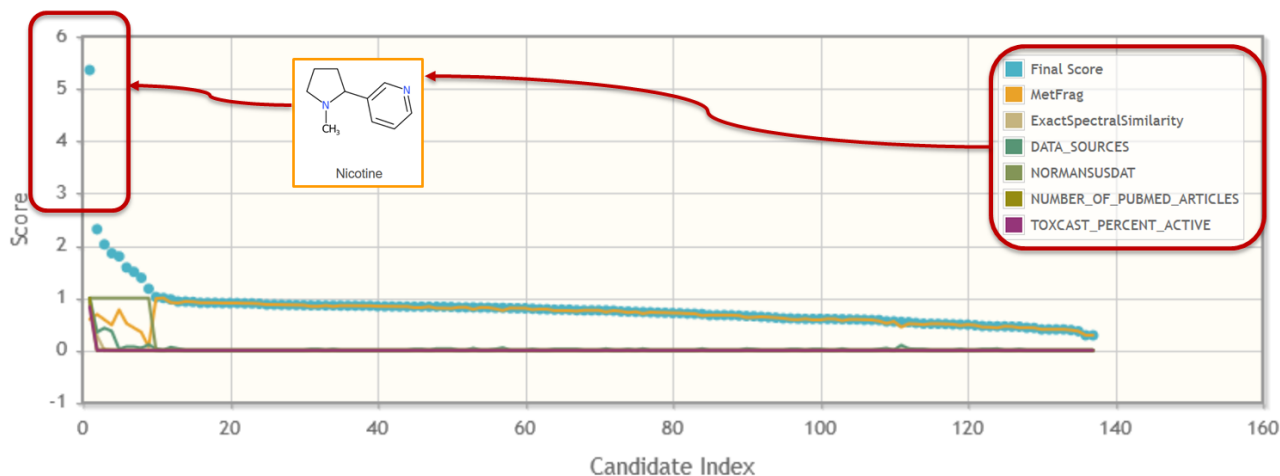
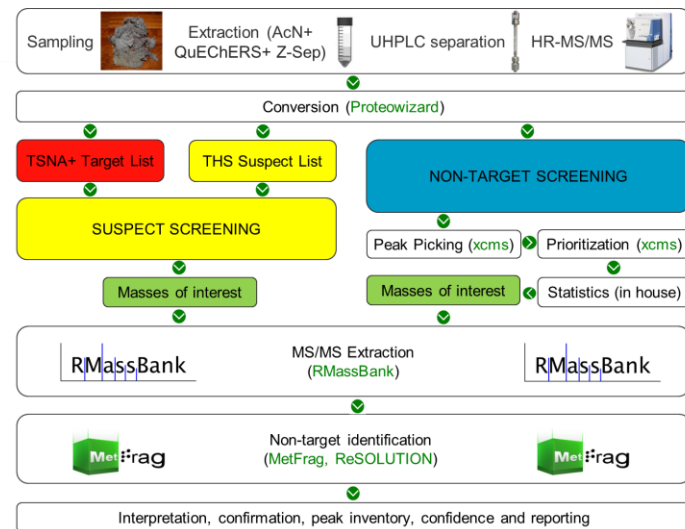
Interpretation, confirmation, peak inventory, confidence and reporting





# Perspectives: Identifying Small Molecules in NTS

- Many comprehensive workflows
  - *I have presented just one of them!*
- Annotation of “known unknowns” is now relatively “quick”:
  - *Especially with well-chosen suspect screening and metadata*



- The bottleneck is still in expt. design and interpretation
  - *But in the meantime we can do some pretty neat things!*

“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

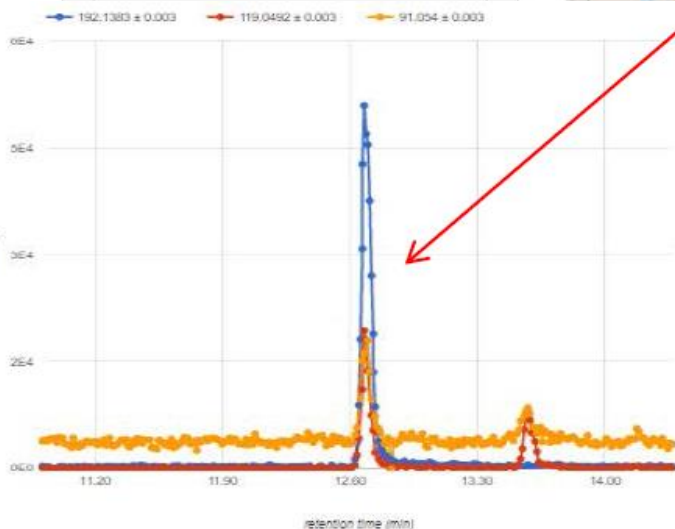
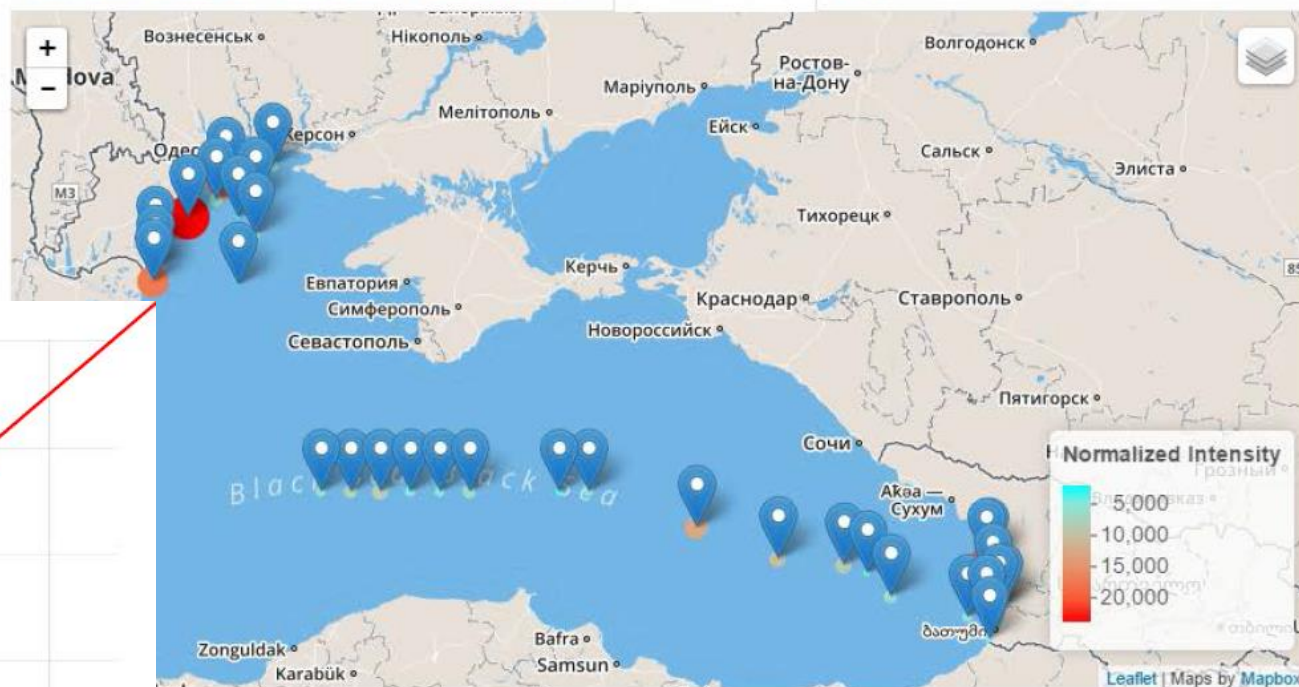
Interactive Map

Help

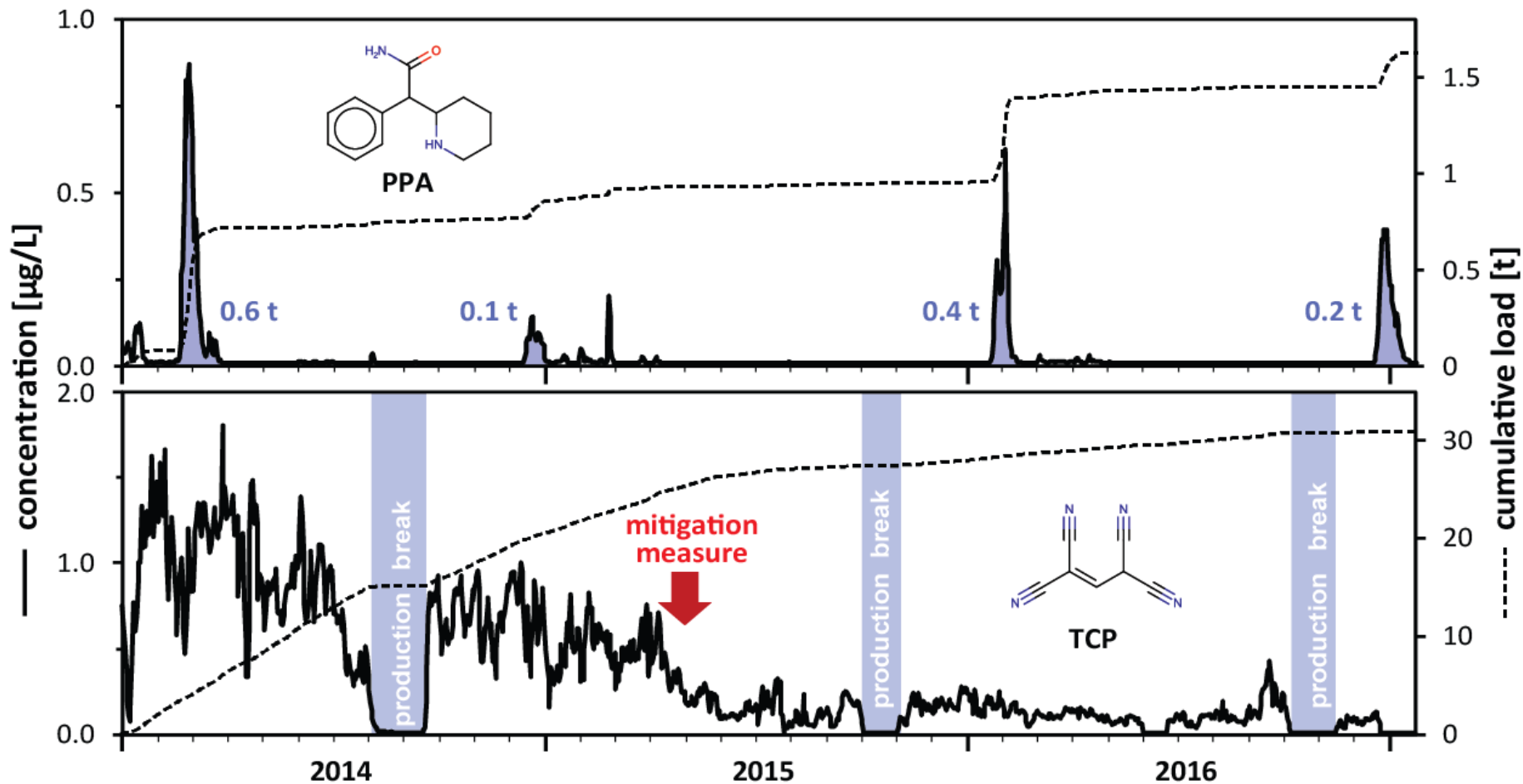
Choose Emerging Substance or input mass of interest and experimental RTI

Substance name or CAS or StdInChIKey

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



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



# New MetaData: Disease-Specific Reference Counts

[https://comptox.epa.gov/dashboard/chemical\\_lists/litminedneuro](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	<a href="#">DTXSID4024983</a>	1032	20	47	140	13	0	4	1	1	0	1	2	4
Ethanol	64-17-5	<a href="#">DTXSID9020584</a>	768	100	23	11	18	26	1	3	20	6	17	54	2
Lead	7439-92-1	<a href="#">DTXSID2024161</a>	740	28	107	68	102	4	2	2	1	3	4	19	30
Lithium	7439-93-2	<a href="#">DTXSID5036761</a>	689	30	50	9	22	5	36	13	25	6	93	12	15
Valproic Acid	76584-70-8	<a href="#">DTXSID70227388</a>	666	32	10	3	65	6	10	18	45	5	18	4	2
1-Methyl-4-phenylpiperazine	28289-54-5	<a href="#">DTXSID8040933</a>	638	1	24	0	11	0	6	289	0	0	5	0	1
Vincristine	2068-78-2	<a href="#">DTXSID8044331</a>	567	17	59	125	15	5	1	1	5	3	2	1	8
Phenytoin	57-41-0	<a href="#">DTXSID8020541</a>	560	37	24	25	16	9	3	1	9	3	8	4	6
Haloperidol	52-86-8	<a href="#">DTXSID4034150</a>	555	6	6	1	10	6	153	51	4	4	11	1	0
Cocaine	50-36-2	<a href="#">DTXSID2038443</a>	530	151	16	0	8	0	2	3	3	8	6	12	11
Aspirin	50-78-2	<a href="#">DTXSID5020108</a>	489	8	3	0	3	2	2	0	9	4	1	0	5
Paclitaxel	33069-62-4	<a href="#">DTXSID9023413</a>	485	4	43	217	9	14	0	0	0	0	0	1	2
Aluminum	7429-90-5	<a href="#">DTXSID3040273</a>	477	13	41	1	105	4	0	0	1	0	1	13	12
Lidocaine	6108-05-0	<a href="#">DTXSID80209953</a>	464	150	26	15	3	2	0	0	8	4	6	2	10
Methotrexate	59-05-2	<a href="#">DTXSID4020822</a>	451	17	25	1	79	4	0	1	5	0	1	9	18
Mercury	7439-97-6	<a href="#">DTXSID1024172</a>	450	6	79	22	23	2	3	5	2	2	38	7	25



▼ PubChem Compound TOC	?	33,765,953
▶ Agrochemical Information	?	2,002
▶ Biologic Description	?	1,539,532
▶ Biological Test Results	?	3,467,416
▶ Biomolecular Interactions and Pathways	?	109,610
▶ Chemical and Physical Properties	?	237,729
▶ Classification	?	18,569,627
▶ Diseases	?	
▶ Drug and Medication Information	?	15,955
▶ Food Additives and Ingredients	?	7,447
▶ Identification	?	5,746
▶ Information Sources	?	20,654,780
▶ Literature	?	1,668,437
▶ Names and Identifiers	?	1,310,169
▶ Patents	?	22,144,888
▶ Pharmacology and Biochemistry	?	130,367
▶ Related Records	?	5,297,096
▶ Safety and Hazards	?	125,607
▶ Spectral Information	?	761,478
▶ Structures	?	5,926,225
▶ Toxicity	?	114,554
▶ Use and Manufacturing	?	108,745
Chemical Safety	?	122,739

- 102 million ...
  - OR
  - Most relevant/annotated
- => PubChemLite

tier0: 316 K

tier1: 360 K



January 14, 2020

Dataset Open Access

## PubChemLite tier0 and tier1

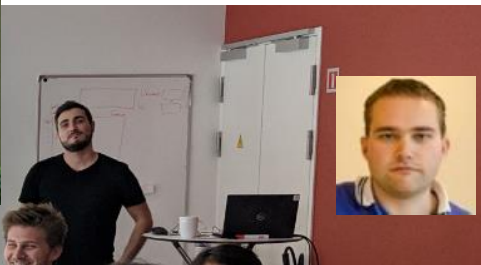
Bolton, Evan; Schymanski, Emma

PubChemLite is a subset of PubChem (<https://pubchem.ncbi.nlm.nih.gov/>) selected from major categories of the Table of Contents page at the PubChem Classification Browser (<https://pubchem.ncbi.nlm.nih.gov/classification/#hid=72>). So far we are providing two "flavours":

tier0 is 316,810 compounds (14 Jan 2020) compiled from 7 categories: AgroChemInfo, DrugMedicInfo, FoodRelated, PharmacolInfo, SafetyInfo, ToxicityInfo, KnownUse

tier1 is 363,911 compounds (14 Jan 2020) compiled from 8 categories (tier0 + BioPathway): AgroChemInfo, BioPathway, DrugMedicInfo, FoodRelated, PharmacolInfo, SafetyInfo, ToxicityInfo, KnownUse

PubChemCIDs have been collapsed by InChIkey first block, reporting the structure from the most annotated CID, plus related CIDs. Entries that will be ignored by MetFrag (salts, disconnected substances) or cause errors (e.g. transition metals) have been removed. The Patent and PubMed ID counts are extracted from files on the PubChem FTP site. The "AnnoTypeCount" term counts how many of the categories are represented, the subsequent column (named per category) counts the number of annotation categories available in the next sub-category of the TOC entry.



ECI@LCSB



Luxembourg National Research Fund



U.S. National Library of Medicine  
National Center for Biotechnology Information



[emma.schymanski@uni.lu](mailto:emma.schymanski@uni.lu)

@ESchymanski @soniatorres @noeliaramz

Further Information:

DOI: [10.5281/zenodo.3613472](https://doi.org/10.5281/zenodo.3613472)

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

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

[https://wwen.uni.lu/lcsb/research/environmental\\_cheminformatics](https://wwen.uni.lu/lcsb/research/environmental_cheminformatics)



UNIVERSITAT ROVIRA I VIRGILI



IISPV  
INSTITUT D'INVESTIGACIÓ SANITÀRIA PERE VIRGILI



UNIVERSITÉ DU LUXEMBOURG



# Community Efforts!



# Supporting Info

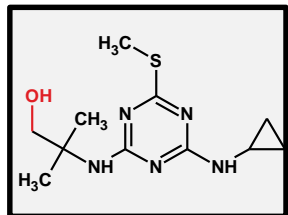
---



Example

Identification confidence

Minimum data requirements



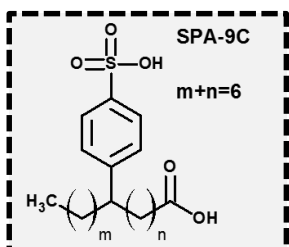
**Level 1: Confirmed structure**  
by reference standard

MS, MS<sup>2</sup>, RT, Reference Std.

**Level 2: Probable structure**

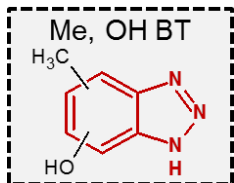
a) by library spectrum match  
b) by diagnostic evidence

MS, MS<sup>2</sup>, Library MS<sup>2</sup>  
MS, MS<sup>2</sup>, Exp. data



**Level 3: Tentative candidate(s)**  
structure, substituent, class

MS, MS<sup>2</sup>, Exp. data



**Level 4: Unequivocal molecular formula**

MS isotope/adduct

C<sub>6</sub>H<sub>5</sub>N<sub>3</sub>O<sub>4</sub>

**Level 5: Exact mass of interest**

MS

192.0757

# "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 <sub>10</sub> H <sub>14</sub> N <sub>2</sub>		4.3349	Peaks: 18 / 23 Fragments Scores Download
2	 Phenylpiperazine	DTXSID40176612 DTXSID40193102 DTXSID90216632 DTXSID50291046 DTXSID00293111 DTXSID50296613  InChIKeyBlock1 = YZTJYBJCZXZGCT	162.11576	C <sub>12</sub> H <sub>16</sub> N <sub>2</sub>			
3	 N'-(2,4-Dimethylphenyl)-N-methylformamide	DTXSID1037696 DTXSID10199510  InChIKeyBlock1 = JIIOLEGNRQDIP	162.11576	C <sub>12</sub> H <sub>16</sub> N <sub>2</sub>			

**LEGEND:** Name, SMILES  
DTXSID | InChIKey 1<sup>st</sup> Block  
CAS | Monois. Mass | logP | Sources  
Data on: Toxicity | Exposure | Bioassays

**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

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

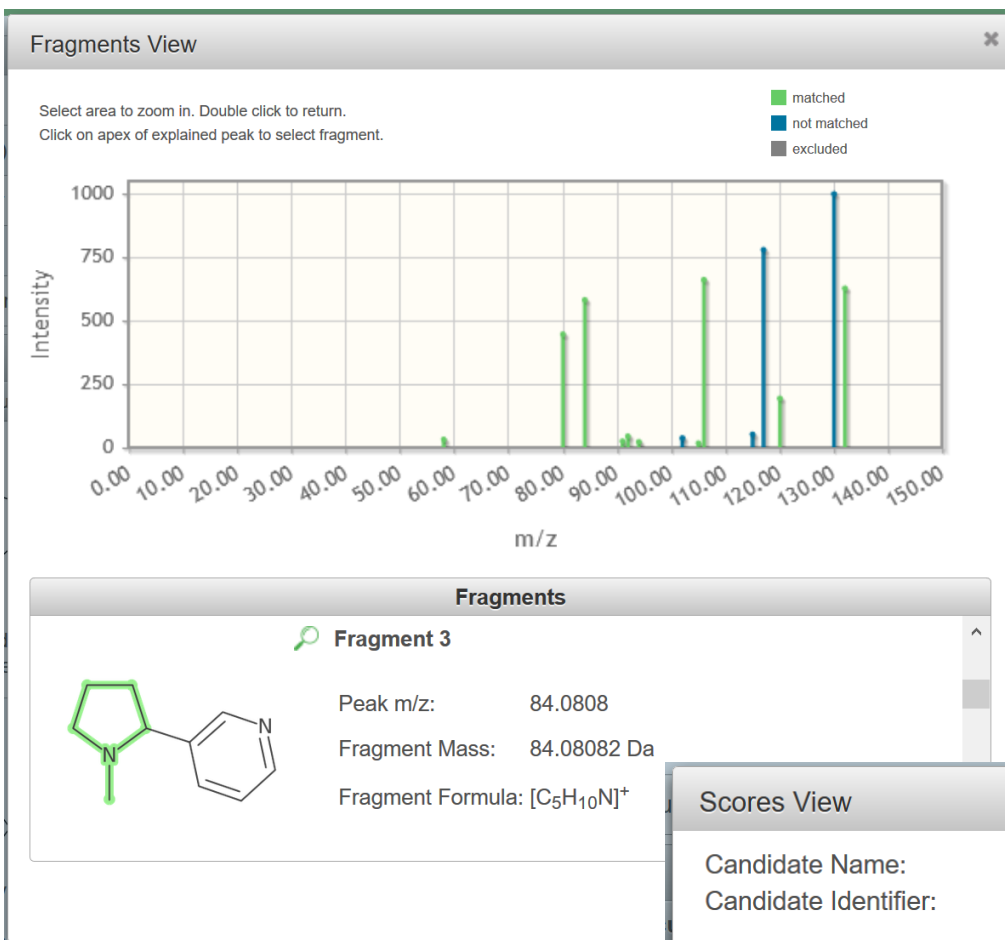
**MS-ready DL-Nicotine**  
CN1CCCC1C1=CN=CC=C1  
 DTXSID3048154 | SNICXCGAKADSCV  
 22083-74-5 | **162.1157** | 0.953 | **9**  
 Tox: yes | Expo: no | Bioassay: yes

**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

**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

**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

# MetFrag – Example with Nicotine VI



### Scores View

Candidate Name: 3-(1-methylpyrrolidin-2-yl)pyridine  
Candidate Identifier: 942

	Name	Normalized Value	Raw Value
🔍	MetFrag	0.5197	332.8153
🔍	SpectralSimilarity	0.9712	6.4214
🔍	ExactSpectralSimilarity	0.9284	0.9284
🔍	PatentsCount	0.9991	71329.0
🔍	PubMedReferenceCount	1.0	18271.0

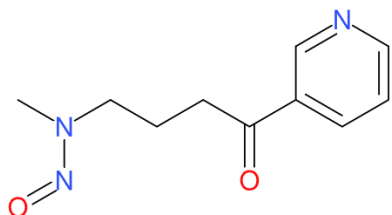
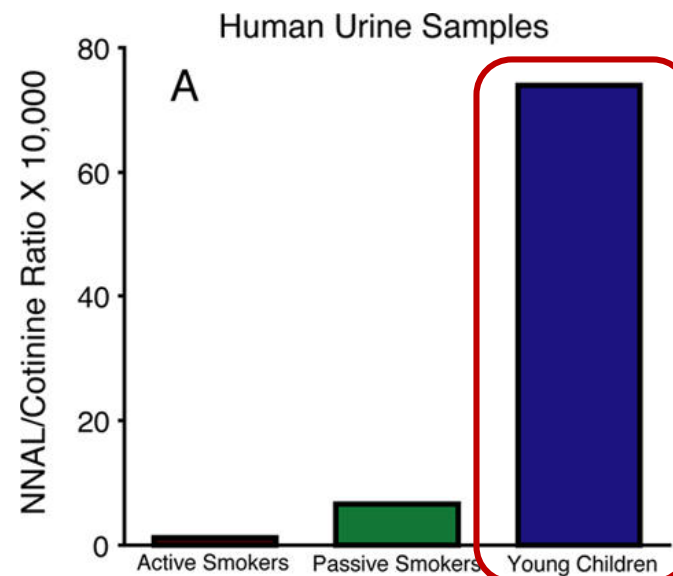
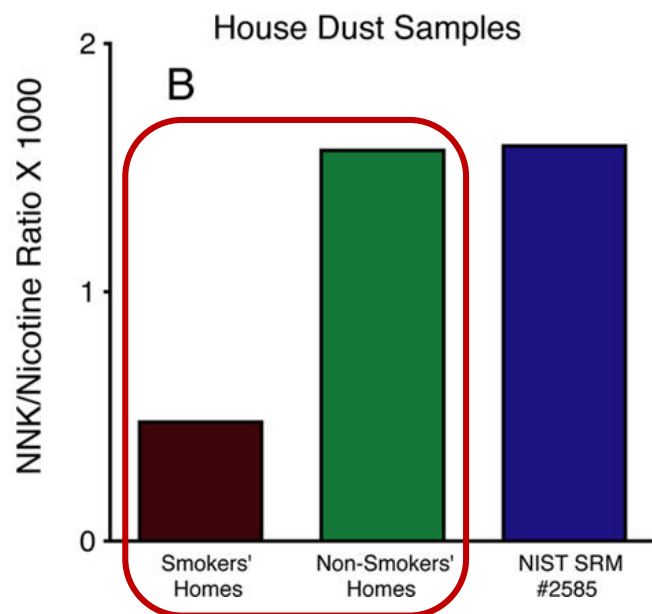
# Why THS is different to SHS



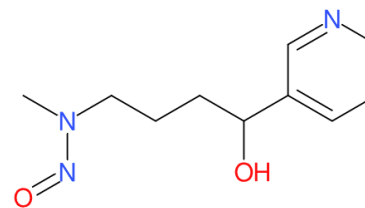
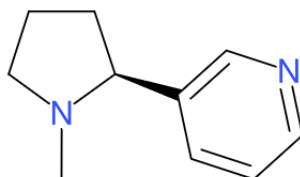
**Ratio  
NNK : Nicotine  
in house dust**



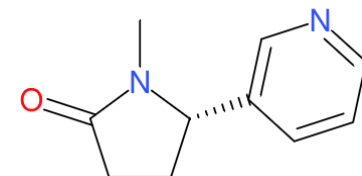
**Ratio of  
NNAL : Cotinine  
in urine**



**NNK : Nicotine**



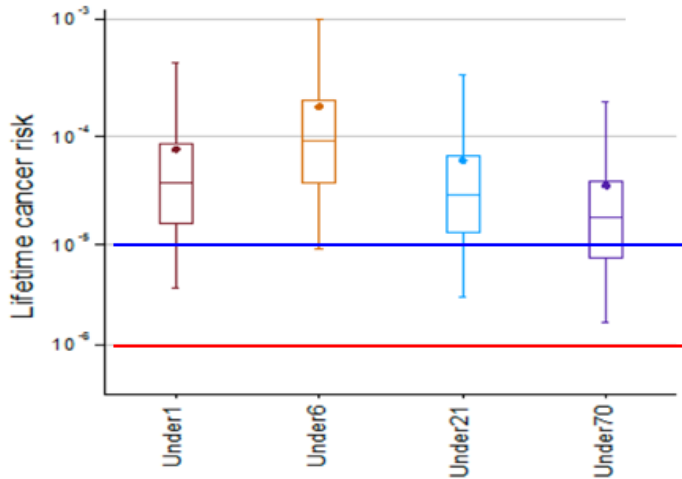
**NNAL : Cotinine**



## Nicotine exposure (ng/kg-day)



Age (years old)	Mean	Max.	Mean	Max.
< 1	129	1637	11	25
1-5	136	1729	12	27
6-21	83	1048	7	16
22-70	80	1030	7	16

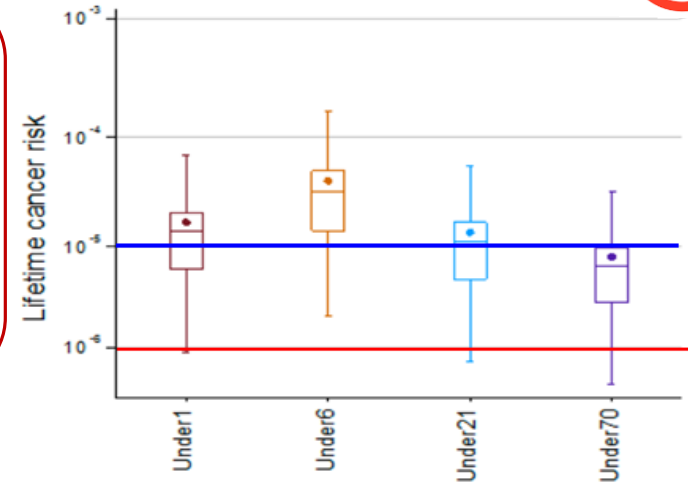


Lifetime cancer  
risk by age group.

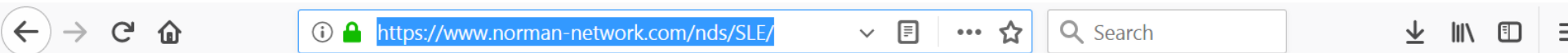
Threshold values

WHO:  $10^{-5}$

USEPA:  $10^{-6}$



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



[NORMAN WEBSITE](#) | 
 [NORMAN DATABASE SYSTEM](#) | 
 [HOME](#) | 
 [LOGIN](#)

NORMAN SUBSTANCE DATABASE

## NORMAN Suspect List Exchange – NORMAN SLE

The NORMAN Suspect List Exchange ([NORMAN-SLE](#)) was established in 2015 as a central access point for NORMAN members (and others) to find suspect lists relevant for their environmental monitoring question. This Exchange documents all individual collections that (will) form a part of [NORMAN SusDat](#), the merged [NORMAN Substance Database](#) (DOI: [10.5281/zenodo.2664077](#)).

If you have any feedback or a list that you would like to have included, please contact [suspects@normandata.eu](mailto:suspects@normandata.eu)

No.	Abbreviation	Description	Link to full list	Link to InChIKey list	References
S0	SUSDAT	<b>Merged NORMAN Suspect List: SusDat</b>	<a href="#">Interactive Data table</a> (updating...) CompTox <a href="#">SUSDAT List</a>	<a href="#">MS-ready InChIKeys</a> (1/03/2018)	A merged list of >40,000 structures from suspect lists. See <a href="#">interactive version</a> . 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: <a href="#">10.5281/zenodo.2664077</a>
S52	THSMOKE	<b>Thirdhand Smoke (THS) Compounds</b>	THSMOKE <a href="#">XLSX</a> , <a href="#">CSV</a> (06/05/2019) CompTox <a href="#">THSMOKE List</a>	THSMOKE <a href="#">InChIKeys</a> (06/05/2019)	Thirdhand Smoke (THS, the tobacco-related gases and particles that become embedded in materials), suspect list compiled by Sonia Torres and Noelia Ramirez (IISPV-URV) and Emma Schymanski (LCSB).  DOI: <a href="#">10.5281/zenodo.2669466</a>

# THSMOKE on CompTox Chemicals Dashboard

[https://comptox.epa.gov/dashboard/chemical\\_lists/THSMOKE](https://comptox.epa.gov/dashboard/chemical_lists/THSMOKE)

## NORMAN: Thirdhand Smoke (THS) Compounds: Suspect List

### List Details

**Description:** A collection of compounds for mass spectrometry suspect screening of Thirdhand Smoke (THS, the tobacco-related gases and particles that become embedded in materials), compiled by S. Torres, N. Ramirez, Institut D'investigacio Sanitaria Pere Virgili at Universitat Rovira i Virgili (IISPV-URV) and E. Schymanski (Luxembourg Center for Systems Biomedicine, LCSB)

**Number of Chemicals:** 95

Select all

Download

Send to Batch Search

Default

95 chemicals

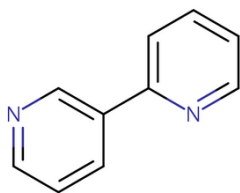
DTXSID

Mass

Sources

Hide chemicals that are:

Filter by Name

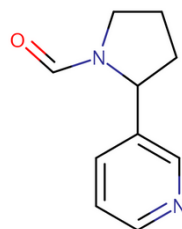


2,3'-Bipyridine

DTXSID:DTXSID00206823

Mass:156.068748

Sources:16

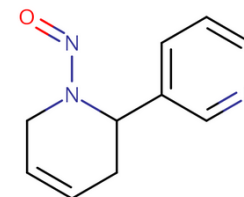


Nornicotine, N-formyl

DTXSID:DTXSID30336006

Mass:176.094963

Sources:3



1-Nitroso-1,2,3,6-tetrahydro-2,3'-bipyrid...

DTXSID:DTXSID40868005

Mass:189.090212

Sources:4