http://www.orcid.org/0000-0002-2668-4821

EPA
United States
Environmental Protection
Agency

# Non-targeted analysis supported by data and cheminformatics delivered via the US EPA CompTox Chemicals Dashboard

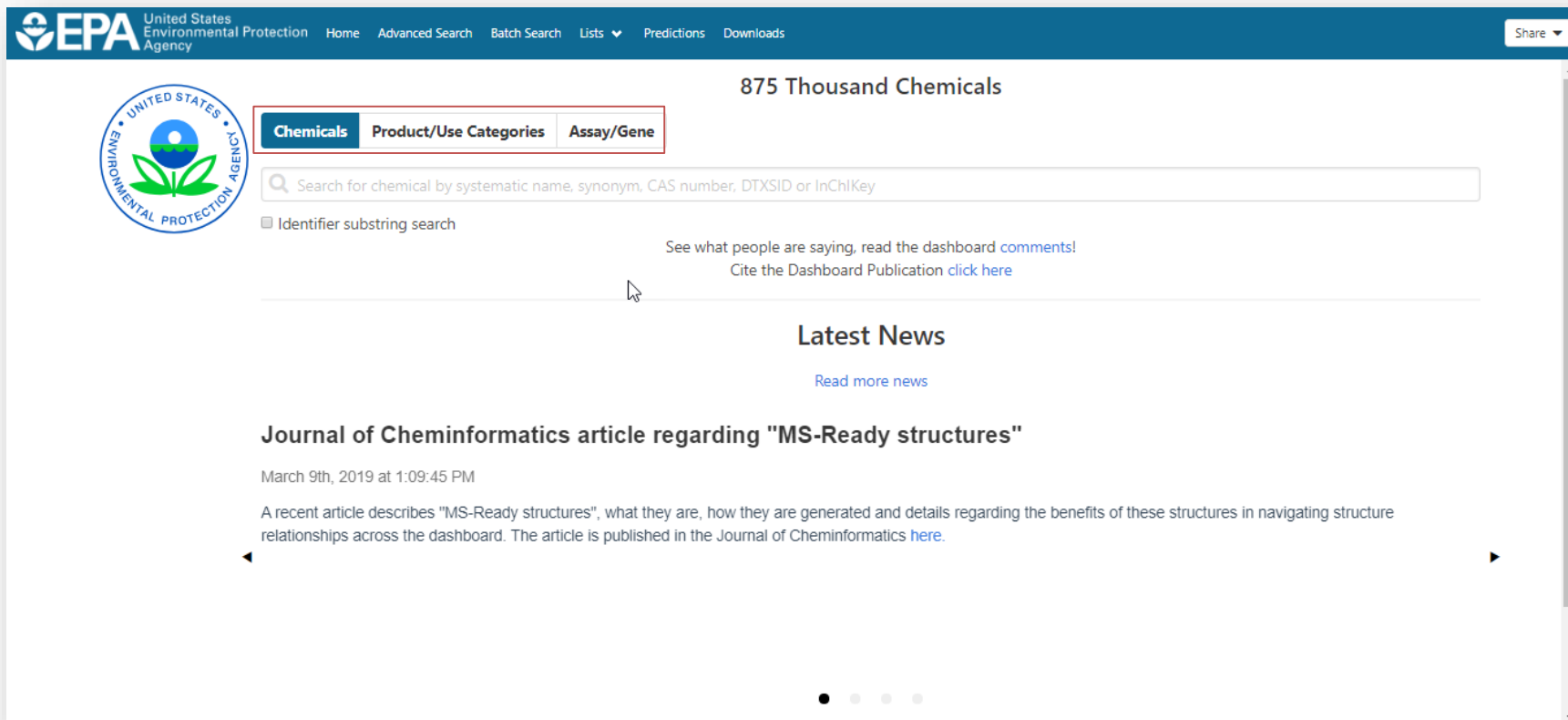**Antony Williams**, *Alex Chao, Tom Transue, Tommy Cathey, Elin Ulrich and Jon Sobus*

1) *National Center for Computational Toxicology, U.S. Environmental Protection Agency, RTP, NC*
2) *Oak Ridge Institute of Science and Education (ORISE) Research Participant, RTP, NC*
3) *GDIT, Research Triangle Park, North Carolina, United State*
4) *National Exposure Research Laboratory, U.S. Environmental Protection Agency, RTP, NC*

*August 2019*
*ACS Fall Meeting, San Diego*

# An intro to the Dashboard

- Freely available web-based database from the National Center for Computational Toxicology

- Providing data for 875,000 substances including
  - Experimental and predicted physicochemical properties
  - *In vivo* toxicity data harvested from dozens of public resources
  - *In vitro* bioactivity data for thousands of chemicals and assays
  - Exposure data including chemicals in consumer products
  - Real time predictions for >20 physchem and toxicological endpoints

- Dashboard is used by mass spectrometrists for chemical identification

- A quick view of general capabilities…

# CompTox Chemicals Dashboard
https://comptox.epa.gov/dashboard

875k Chemical Substances

# Detailed Chemical Pages

# Access to Chemical Hazard Data

# Sources of Exposure to Chemicals

# Link Access



**Links based on chemical identifiers to dozens of online resources – including analytical data**

# MassBank of North America
# https://mona.fiehnlab.ucdavis.edu

# "MS-ready" structures

Journal of Cheminformatics

**METHODOLOGY**    **Open Access**

## "MS-Ready" structures for non-targeted high-resolution mass spectrometry screening studies

Andrew D. McEachran[1,2*], Kamel Mansouri[1,2,3], Chris Grulke[2], Emma L. Schymanski[4], Christoph Ruttkies[5] and Antony J. Williams[2*]

8

# Overview of MS-Ready Structures

- All structure-based chemical substances are algorithmically processed to
  - Split multicomponent chemicals into individual structures
  - Desalt and neutralize individual structures
  - Remove stereochemical bonds from all chemicals

- MS-Ready structures are then mapped to original substances to provide a path between chemicals detected by mass spectrometry to original substances

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
DTXSID004635| SNICXCGAKADSCV
25162-00-9 | **162.1157**| 0.929| **20**
Tox: **no**| Expo: **yes**| Bioassay: **yes**

HCl
Nicotine hydrochloride
Cl.CN1CCC[C@H]1C1=CN=CC=C1
DTXSID602093| HDJBTCAJIMNXEW
2820-51-1 | **198.0924** | 0.929| **9**
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**

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 Mappings from Details Page

# Two MS-Ready Mappings Set

# MS-Ready Mappings Set
## All substances containing component

# *Mass/Formula Searching and Metadata Ranking*

# Advanced Searches
## **Mass** Search

# Advanced Searches
## **Mass** Search

# MS-Ready Structures for
**Formula Search**

# MS-Ready Mappings

- **EXACT Formula**: C10H16N2O8: **3** Hits

# MS-Ready Mappings

- **Same** Input Formula: C10H16N2O8

- <span style="color:red">**MS Ready Formula**</span> Search: **125** Chemicals

- Exact Formula – 3 hits

- MS-Ready Formula – 125 hits!!
  - ONLY 8 of the 125 are **single component** chemicals
  - 3 are **neutral compounds** and 2 are **charged**

- How can we rank the candidates list?

# *Candidate ranking using metadata*

**RESEARCH ARTICLE**

**Identification of "Known Unknowns" Utilizing Accurate Mass Data and ChemSpider**

# Data Source Ranking of "*known unknowns*"

- A mass and/or formula search is for an ***unknown*** chemical but it is a ***known*** chemical contained within a reference database

- **Most likely** candidate chemicals have the **most** associated data sources, **most** associated literature articles or both

$C_{14}H_{22}N_2O_3$
266.16304

↓

**Chemical Reference Database**

↓

**Sorted candidate structures**

# The original ChemSpider work

| Compound class | Number in class | Average rank | Number of compounds in each position rank-ordered | | | | |
|---|---|---|---|---|---|---|---|
| | | | #1 | #2 | #3 | #4 | #5+ |
| Pharmaceutical drug | 72 | 1.4 | 55 | 9 | 6 | 2 | |
| Industrial chemicals | 42 | 5.5 | 28 | 6 | 3 | | 5 |
| Personal care products | 8 | 6.1 | 3 | 1 | | | 4 |
| Steroid hormones | 7 | 1.0 | 7 | | | | |
| Perfluorochemicals | 6 | 1.2 | 5 | 1 | | | |
| Pesticides | 12 | 2.3 | 6 | 2 | 3 | | 1 |
| Veterinary drugs | 3 | 1.3 | 2 | 1 | | | |
| Dyes | 2 | 1.0 | 2 | | | | |
| Food product/natural compounds | 4 | 3.8 | 2 | | | 1 | 1 |
| Illicit drugs | 2 | 2.0 | 1 | | 1 | | |
| Misc. molecules | 3 [a] | 1.3 | 2 | 1 | | | |

# Is a bigger database better?

- ChemSpider was 26 million chemicals for the original work

- Much BIGGER today

- Is bigger better??

- Are there other metadata to use for ranking?

**71Million**
chemical structures

# Using Metadata for Ranking

- Chosen dashboard metadata to rank candidates
  - Associated data sources
    - Lists in the underlying database (more about lists later)
    - Associated data sources in PubChem
    - Specific source types (e.g. water, surfactants, pesticides)
  - Number of associated literature articles (Pubmed)
  - **Chemicals in the environment** – the number of products/categories containing the chemical is an important source of data (from CPDat database)

# Identification ranks for **1783 chemicals** using multiple data streams

**Data Sources alone rank ~75% of the chemicals as Top Hit**



**DS: Data Sources**
**PC: PubChem**
**PM: PubMed**
**STOFF: DB**
**KEMI: DB**

# *Comparing* Search Performance

Anal Bioanal Chem (2017) 409:1729–1735
DOI 10.1007/s00216-016-0139-z

CrossMark

**RAPID COMMUNICATION**

# Identifying known unknowns using the US EPA's CompTox Chemistry Dashboard

Andrew D. McEachran[1] · Jon R. Sobus[2] · Antony J. Williams[3]

- When dashboard contained 720k chemicals
- Only **3%** of ChemSpider size
- What was the comparison in performance?

# SAME dataset for comparison

| Compound class | Number in class | Average rank | Number of compounds in each position rank-ordered | | | | |
|---|---|---|---|---|---|---|---|
| | | | #1 | #2 | #3 | #4 | #5+ |
| Pharmaceutical drug | 72 | 1.4 | 55 | 9 | 6 | 2 | |
| Industrial chemicals | 42 | 5.5 | 28 | 6 | 3 | | 5 |
| Personal care products | 8 | 6.1 | 3 | 1 | | | 4 |
| Steroid hormones | 7 | 1.0 | 7 | | | | |
| Perfluorochemicals | 6 | 1.2 | 5 | 1 | | | |
| Pesticides | 12 | 2.3 | 6 | 2 | 3 | | 1 |
| Veterinary drugs | 3 | 1.3 | 2 | 1 | | | |
| Dyes | 2 | 1.0 | 2 | | | | |
| Food product/natural compounds | 4 | 3.8 | 2 | | | 1 | 1 |
| Illicit drugs | 2 | 2.0 | 1 | | 1 | | |
| Misc. molecules | 3 [a] | 1.3 | 2 | 1 | | | |

EXACTLY THE SAME DATASET

28

# How did performance compare?

| | Mass-based searching | | Formula-based searching | |
|---|---|---|---|---|
| | Dashboard | ChemSpider | Dashboard | ChemSpider |
| Average rank position | 1.3 | 2.2[a] | 1.2 | 1.4 |
| Percent in #1 position | 85% | 70% | 88% | 80% |

[a] Average rank in ChemSpider shown here does not include an outlier where the rank was 201, when added the average rank position is 3.5

**For the same 162 chemicals, Dashboard outperforms ChemSpider for both Mass and Formula Ranking**

# How did performance compare?

Summary statistics and rank-ordered position in the CompTox Chemistry Dashboard and ChemSpider of the 89 compound subset from the Little et al. [7] study

| | | Average rank | Number in each position rank-ordered | | | | |
|---|---|---|---|---|---|---|---|
| | | (±SD) | #1 | #2 | #3 | #4 | #5+ |
| Mass-based | Dashboard | $1.2 \pm 0.7$ | $77^a$ | 5 | 3 | 3 | |
| | ChemSpider | $2.2 \pm 6.1^b$ | 68 | 8 | 7 | 1 | 5 |
| Formula-based | Dashboard | $1.1 \pm 0.4$ | $78^a$ | 8 | 2 | | |
| | ChemSpider | $1.3 \pm 1.0$ | 77 | 8 | 2 | 1 | 2 |

[a]One chemical (tephrosin) not present in the Dashboard

# Data Quality is important

- ## Data quality in free web-based databases!

## Machines first, humans second: on the importance of algorithmic interpretation of open chemistry data

Alex M Clark ✉, Antony J Williams and Sean Ekins

31

# Public Databases require curation

- There is significant bloating in the public databases because of lack of curation

- The number of hits retrieved based on mass or formula searching can explode based on poorly represented chemicals – especially stereochemistry issues

- MS-Ready structures will map back to multiple versions of "the same chemical".

# Will the correct Microcystin LR Stand Up?
# ChemSpider Skeleton Search

# Comparing ChemSpider Structures



| ChemSpiderID | Standard InChIKey Stereolayer |
|---|---|
| WIKIPEDIA | t28-,29-,30-,31+,34-,35-,36+,37+,38-,40+ |
| CompTox | t28-,29-,30-,31+,34-,35-,36+,37+,38-,40+ |
| 4941647 | t28-,29-,30-,31+,34-,35-,36+,37+,38-,40+ |
| 393078 | t28-,29-,30-,31+,34-,35-,36+,**37-**,38-,40+ |
| 57618348 | t28-,29-,30-,31+,34-,35-,36+,**37-**,38-,40+ |
| 29342071 | t28-,29-,30-,31+,**34+**,35-,36+,**37-**,38-,40+ |
| 7987594 | t28-,**29?,30?**,31+,**34?**,35-,**36?,37-**,38-,**40?** |
| 22900854 | t28-,**29?,30+,31-,34+,35+,36-,37-**,38-,**40-** |
| 19692240 | NONE |
| 2831283 | NONE |

# Comparing ChemSpider Structures

| ChemSpiderID | InChIKey | # Stereocenters | # Different |
|---|---|---|---|
| **WIKIPEDIA** | ZYZCGGRZINLQBL-JCGNTXOTSA-N | 10/10 | 0 |
| **CompTox** | ZYZCGGRZINLQBL-JCGNTXOTSA-N | 10/10 | 0 |
| 4941647 | ZYZCGGRZINLQBL-JCGNTXOTSA-N | 10/10 | 0 |
| 393078 | ZYZCGGRZINLQBL-GWRQVWKTSA-N | 10/10 | 1 |
| 57618348 | ZYZCGGRZINLQBL-UPPCHHEJSA-N | 10/10 | 1 |
| 29342071 | ZYZCGGRZINLQBL-IIJTUTQBSA-N | 10/10 | 2 |
| 7987594 | ZYZCGGRZINLQBL-BESLYTPASA-N | 5/10 | 6 |
| 22900854 | ZYZCGGRZINLQBL-QAXSDTKVSA-N | 9/10 | 8 |
| 19692240 | ZYZCGGRZINLQBL-ORZJCNCZSA-N | 0/10 | 10 |
| 2831283 | ZYZCGGRZINLQBL-UHFFFAOYSA-N | 0/10 | 10 |

# Other Searches

# *Batch Searching mass and formula*

# Batch Searching

- Singleton searches are useful but we work with **thousands** of masses and formulae!

- Typical questions
  - What is the list of chemicals for the formula $C_xH_yO_z$
  - What is the list of chemicals for a mass +/- error
  - Can I get chemical lists in Excel files? In SDF files?
  - Can I include properties in the download file?

# Batch Searching Formula/Mass

# Searching batches using MS-Ready Formula (or mass) searching

| | A | B | C | D | E | F | G |
|---|---|---|---|---|---|---|---|
| 1 | INPUT | DTXSID | CASRN | PREFERRED NAME | MOL FORMULA | MONOISOTOPIC MASS | DATA SOURCES |
| 2 | C14H22N2O3 | DTXSID2022628 | 29122-68-7 | Atenolol | C14H22N2O3 | 266.163042576 | 46 |
| 3 | C14H22N2O3 | DTXSID0021179 | 6673-35-4 | Practolol | C14H22N2O3 | 266.163042576 | 32 |
| 4 | C14H22N2O3 | DTXSID4048854 | 841-73-6 | Bucolome | C14H22N2O3 | 266.163042576 | 20 |
| 5 | C14H22N2O3 | DTXSID1045407 | 13171-25-0 | Trimetazidine dihydrochloride | C14H24Cl2N2O3 | 338.116398 | 19 |
| 6 | C14H22N2O3 | DTXSID0045753 | 56715-13-0 | R-(+)-Atenolol | C14H22N2O3 | 266.163042576 | 19 |
| 7 | C14H22N2O3 | DTXSID2048531 | 5011-34-7 | Trimetazidine | C14H22N2O3 | 266.163042576 | 14 |
| 8 | C14H22N2O3 | DTXSID10239405 | 93379-54-5 | Esatenolol | C14H22N2O3 | 266.163042576 | 12 |
| 9 | C14H22N2O3 | DTXSID50200634 | 52662-27-8 | N-(2-Diethylaminoethyl)-2-(4-hydroxyphenoxy)acetamide | C14H22N2O3 | 266.163042576 | 7 |
| 10 | C14H22N2O3 | DTXSID4020111 | 51706-40-2 | dl-Atenolol hydrochloride | C14H23ClN2O3 | 302.1397203 | 6 |
| 11 | C14H22N2O3 | DTXSID1068693 | 51963-82-7 | Benzenamine, 2,5-diethoxy-4-(4-morpholinyl)- | C14H22N2O3 | 266.163042576 | 5 |
| 12 | C18H34N2O6S | DTXSID3023215 | 154-21-2 | Lincomycin | C18H34N2O6S | 406.213757997 | 35 |
| 13 | C18H34N2O6S | DTXSID7047803 | 859-18-7 | Lincomycin hydrochloride | C18H35ClN2O6S | 442.1904357 | 22 |
| 14 | C18H34N2O6S | DTXSID20849438 | 1398534-62-7 | PUBCHEM_71432748 | C18H35ClN2O6S | 442.1904357 | 1 |
| 15 | C10H12N2O | DTXSID1047576 | 486-56-6 | Cotinine | C10H12N2O | 176.094963014 | 40 |
| 16 | C10H12N2O | DTXSID8075330 | 50-67-9 | Serotonin | C10H12N2O | 176.094963014 | 22 |
| 17 | C10H12N2O | DTXSID8044412 | 2654-57-1 | 4-Methyl-1-phenylpyrazolidin-3-one | C10H12N2O | 176.094963014 | 18 |
| 18 | C10H12N2O | DTXSID80165186 | 153-98-0 | Serotonin hydrochloride | C10H13ClN2O | 212.0716407 | 11 |
| 19 | C10H12N2O | DTXSID2048870 | 29493-77-4 | (4R,5S)-4-methyl-5-phenyl-4,5-dihydro-1,3-oxazol-2-amine | C10H12N2O | 176.094963014 | 10 |
| 20 | C10H12N2O | DTXSID10196105 | 443-31-2 | 6-Hydroxytryptamine | C10H12N2O | 176.094963014 | 9 |
| 21 | C10H12N2O | DTXSID90185693 | 31822-84-1 | 1,4,5,6-Tetrahydro-5-phenoxypyrimidine | C10H12N2O | 176.094963014 | 7 |
| 22 | C10H12N2O | DTXSID40178777 | 2403-66-9 | 2-Benzimidazolepropanol | C10H12N2O | 176.094963014 | 7 |
| 23 | C10H12N2O | DTXSID80157026 | 13140-86-8 | N-Cyclopropyl-N'-phenylurea | C10H12N2O | 176.094963014 | 6 |
| 24 | C10H12N2O | DTXSID30205607 | 570-14-9 | 4-Hydroxytryptamine | C10H12N2O | 176.094963014 | 6 |
| 25 | C14H18N4O3 | DTXSID5023900 | 17804-35-2 | Benomyl | C14H18N4O3 | 290.137890456 | 68 |
| 26 | C14H18N4O3 | DTXSID3023712 | 738-70-5 | Trimethoprim | C14H18N4O3 | 290.137890456 | 51 |
| 27 | C14H18N4O3 | DTXSID40209671 | 60834-30-2 | Trimethoprim hydrochloride | C14H19ClN4O3 | 326.1145682 | 8 |
| 28 | C14H18N4O3 | DTXSID70204210 | 55687-49-5 | Benzenemethanol, 4-((2,4-diamino-5-pyrimidinyl)methyl)-2, | C14H18N4O3 | 290.137890456 | 5 |
| 29 | C14H18N4O3 | DTXSID20152671 | 120075-57-2 | 6-Methoxy-4-(3-(N,N-dimethylamino)propylamino)-5,8-quina | C14H18N4O3 | 290.137890456 | 4 |
| 30 | C14H18N4O3 | DTXSID30213742 | 63931-79-3 | 1H-1,2,4-Benzotriazepine-3-carboxylic acid, 4,5-dihydro-4- | C14H18N4O3 | 290.137890456 | 3 |
| 31 | C14H18N4O3 | DTXSID30219608 | 69449-07-6 | 2,4-Pyrimidinediamine, 5-((3,4,5-trimethoxyphenyl)methyl)- | C14H20N4O4 | 308.14845514 | 3 |
| 32 | C14H18N4O3 | DTXSID20241155 | 94232-27-6 | L-Aspartic acid, compound with 5-((3,4,5-trimethoxyphenyl | C18H25N5O7 | 423.175398165 | 3 |
| 33 | C14H18N4O3 | DTXSID80241156 | 94232-28-7 | L-Glutamic acid, compound with 5-((3,4,5-trimethoxypheny | C19H27N5O7 | 437.191048229 | 3 |
| 34 | C14H18N4O3 | DTXSID20143781 | 101204-93-7 | 1H-Pyrido(2,3-e)-1,4-diazepine-2,3,5-trione, 4-(2-(diethylam | C14H18N4O3 | 290.137890456 | 3 |
| 35 | C12H11N7 | DTXSID6021373 | 396-01-0 | Triamterene | C12H11N7 | 253.107593382 | 52 |
| 36 | C12H11N7 | DTXSID00204465 | 5587-93-9 | Ampyrimine | C12H11N7 | 253.107593382 | 7 |
| 37 | C12H11N7 | DTXSID5064621 | 7300-26-7 | Benzenamine, 4-azido-N-(4-azidophenyl)- | C12H9N7 | 251.091943318 | 4 |
| 38 | C12H11N7 | DTXSID00848025 | 90293-82-6 | Sulfuric acid--6-phenylpteridine-2,4,7-triamine (1/1) | C12H13N7O4S | 351.074973101 | 1 |
| 39 | C12H11N7 | DTXSID50575293 | 92310-83-3 | (1E)-N-Phenyl-1,2-bis(1H-1,2,4-triazol-1-yl)ethan-1-imine | C12H11N7 | 253.107593382 | 1 |
| 40 | C8H9NO2 | DTXSID2020006 | 103-90-2 | Acetaminophen | C8H9NO2 | 151.063328534 | 75 |
| 41 | C8H9NO2 | DTXSID6025567 | 134-20-3 | Methyl 2-aminobenzoate | C8H9NO2 | 151.063328534 | 50 |

# *Mass Spectrometry Related Searches*

# Find me "related structures"
# Formula-Based Search

# Select Chemicals of Interest

# Find me "related structures" Based on Structure Similarity



**Morphine**
57-27-2 | DTXSID9023336
Searched by Approved Name.

| Wikipedia | ◀ |
| --- | --- |
| **Intrinsic Properties** | ◀ |
| **Structural Identifiers** | ◀ |
| **Linked Substances** | ▼ |

**Same Connectivity:** 3 records (based on first layer of InChI)

**Mixtures, Components and Neutralized Forms:** 13 records (based on QSAR ready mappings and with the compound as a component of a mixture)

**MS-Ready Mappings:** DTXCID60196731: 21 records;

**Similar Compounds:** 178 records (based on Tanimoto coefficient >0.8)

| **Presence in Lists** | ◀ |
| --- | --- |

# Find me "related structures" Based on Structure Similarity

# Find me "related structures"
# Structure Similarity – sort on mass

# *Chemical Lists*

# Chemical Lists

## WATER|EPA; Chemicals associated with hydraulic fracturing

Search EPAHFR Chemicals

☐ Identifier substring search

### List Details ▼

**Description:** Chemicals used in hydraulic fracturing fluids and/or identified in produced water from 2005-2013, corresponding to chemicals listed in Appendix H of EPA's Hydraulic Fracking Drinking Water Assessment Final Report (Dec 2016). Citation: U.S. EPA, Hydraulic Fracturing for Oil and Gas: Impacts from the Hydraulic Fracturing Water Cycle on Drinking Water Resources in the United States (Final Report). U.S. Environmental Protection Agency, Washington, D.C. EPA/600/R-16/236F, 2016. https://www.epa.gov/hfstudy

*Note that Appendix H chemical listings in Tables H-2 and H-4 were mapped to current DSSTox content, which has undergone additional curation since the publication of the original EPA HF Report (Dec 2016). In the few cases where a Chemical Name and CASRN from the original report map to distinct substances (as of Jan 2018), both were included in the current EPAHFR chemical listing for completeness; additionally, 34 previously unmapped chemicals in Table H-5 are now registered in DSSTox (all but 2 assigned CASRN) and, thus, have been added to the current EPAHFR listing.
**Number of Chemicals:** 1640

| Alkylbenzenesulfonate, linear | Ammonium chloride | Diammonium citrate | Ammonium hydroxide |
|---|---|---|---|
| DTXSID: DTXSID3020041 | DTXSID: DTXSID0020078 | DTXSID: DTXSID5020079 | DTXSID: DTXSID4020080 |
| PubChem: 0 | PubChem: 82 | PubChem: 19 | PubChem: 83 |
| CPDAT: 83 | CPDAT: 260 | CPDAT: 18 | CPDAT: 857 |

# PFAS lists of Chemicals



## Select List

⬇ Download ▾    Columns ▾

PFAS    📋 Copy Filtered Lists URL

| List Acronym ⇕ | List Name ⇕ | Last Updated ⇕ | Number of Chemicals ⇕ | List Description |
|---|---|---|---|---|
| EPAPFAS75S1 | PFAS\|EPA: List of 75 Test Samples (Set 1) | 2018-06-29 | 74 | PFAS list corresponds to 75 samples (Set 1) submitted for initial testing screens conducted by EPA researchers in collaboration with researchers at the National Toxicology Program. |
| EPAPFAS75S2 | PFAS\|EPA: List of 75 Test Samples (Set 2) | 2019-02-21 | 75 | PFAS list corresponds to a second set of 75 samples (Set 2) submitted for testing screens conducted by EPA researchers in collaboration with researchers at the National Toxicology Program. |
| EPAPFASCAT | PFAS\|EPA Structure-based Categories | 2018-06-29 | 64 | List of registered DSSTox "category substances" representing PFAS categories created using ChemAxon's Markush structure-based query representations. |
| EPAPFASINSOL | PFAS\|EPA: Chemical Inventory Insoluble in DMSO | 2018-06-29 | 43 | PFAS chemicals included in EPA's expanded ToxCast chemical inventory found to be insoluble in DMSO above 5mM. |
| EPAPFASINV | PFAS\|EPA: ToxCast Chemical Inventory | 2018-06-29 | 430 | PFAS chemicals included in EPA's expanded ToxCast chemical inventory and available for testing. |
| EPAPFASRL | PFAS\|EPA: Cross-Agency Research List | 2017-11-16 | 199 | EPAPFASRL is a manually curated listing of mainly straight-chain and branched PFAS (Per- & Poly-fluorinated alkyl substances) compiled from various internal, literature and public sources by EPA researchers and program office representatives. |
| PFASKEMI | PFAS: List from the Swedish Chemicals Agency (KEMI) Report | 2017-02-09 | 2416 | Perfluorinated substances from a Swedish Chemicals Agency (KEMI) Report on the occurrence and use of highly fluorinated substances. |
| PFASMASTER | PFAS Master List of PFAS Substances | 2018-07-26 | 5061 | PFASMASTER is a consolidated list of PFAS substances spanning and bounded by the below lists of current interest to researchers and regulators worldwide. |
| PFASOECD | PFAS: Listed in OECD Global Database | 2018-05-16 | 4729 | OECD released a New Comprehensive Global Database of Per- and Polyfluoroalkyl Substances, (PFASs) listing more than 4700 new PFAS |
| PFASTRIER | PFAS Community-Compiled List (Trier et al., 2015) | 2017-07-16 | 597 | PFASTRIER community-compiled public listing of PFAS (Trier et al, 2015) |

# *Research in Progress*

# Predicted Mass Spectra

http://cfmid.wishartlab.com/



- MS/MS spectra prediction for ESI+, ESI-, and EI
- Predictions generated and stored for >800,000 structures, to be accessible via Dashboard

# Search Expt. vs. Predicted Spectra

# Search Expt. vs. Predicted Spectra

# Spectral Viewer Comparison

# Prototype Development

# Prototype Development

# API services and Open Data

- Present API and web services available at https://actorws.epa.gov/actorws/ but major redevelopment is underway

- Downloadable data available via the downloads page

# Web Services
## https://actorws.epa.gov/actorws/

- Data in UI, JSON and XML format

https://actorws.epa.gov/actorws/dsstox/v02/msready?identifier=80-05-7
https://actorws.epa.gov/actorws/dsstox/v02/msready.json?identifier=80-05-7
https://actorws.epa.gov/actorws/dsstox/v02/msready.xml?identifier=80-05-7

https://actorws.epa.gov/actorws/dsstox/v02/msready?identifier=DTXCID60513
https://actorws.epa.gov/actorws/dsstox/v02/msready.json?identifier=DTXCID60513
https://actorws.epa.gov/actorws/dsstox/v02/msready.xml?identifier=DTXCID60513

https://actorws.epa.gov/actorws/dsstox/v02/msready?identifier=UVOFGKIRTCCNKG-UHFFFAOYSA-N
https://actorws.epa.gov/actorws/dsstox/v02/msready.json?identifier=UVOFGKIRTCCNKG-UHFFFAOYSA-N
https://actorws.epa.gov/actorws/dsstox/v02/msready.xml?identifier=UVOFGKIRTCCNKG-UHFFFAOYSA-N

# InChIKey to DTXCIDs

**https://actorws.epa.gov/actorws/dsstox/v02/msready?identifier=UVOFGKIRTCCNKG-UHFFFAOYSA-N**

| Image | DTXCID | Smiles | Image | MsReady DTXCID | MsReady SMILES |
|---|---|---|---|---|---|
|  | DTXCID60513 | C[NH2+]C.CN(C)C([S-])=S |  | DTXCID0023797 | CN(C)C(S)=S |
|  | DTXCID60513 | C[NH2+]C.CN(C)C([S-])=S |  | DTXCID704057 | CNC |

# *Data and Services used by the Community*

# NORMAN Suspect List Exchange
https://www.norman-network.com/?q=node/236

## NORMAN
**Network of reference laboratories, research centres and related organisations for monitoring of emerging environmental substances**
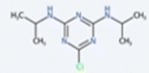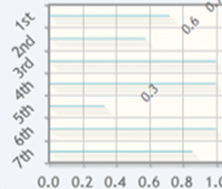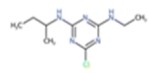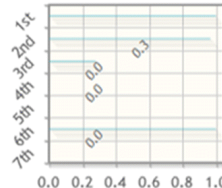
| | | | |
|---|---|---|---|
| **Wastewater Suspect List based on Swedish Product Data** | Wastewater Suspect List **Original File with Mapped DTXSIDs** (12/02/2019) | KEMIWWSUS **InChIKeys** (12/02/2019) | A prioritized list of 1,123 substances relevant for wastewater based on Swedish product registry data, including scores. Provided by Stellan Fischer, KEMI. |
| **Algal toxins list from CompTox** | ALGALTOX **XLSX, CSV** (14/02/2019) CompTox **ALGALTOX List** | ALGALTOX **InChIKeys** (14/02/2019) | List of algal toxins (generated during blooms) from the CompTox Chemicals Dashboard. |
| **CCL 4 Chemical Candidate List** | CCL4 **XLSX, CSV** (14/02/2019) CompTox **CCL4 List** | CCL4 **InChIKeys** (14/02/2019) | Contaminants that are not (yet) regulated in the USA but are known or anticipated to occur in public water systems; from CompTox. |
| **Hydrogen Deuterium Exchange (HDX) Standard Set** | HDXNOEX **XLSX, CSV** (14/02/2019) CompTox **HDXNOEX List** CompTox **HDXEXCH List** | HDXNOEX **InChIKeys** (14/02/2019) | Environmental standard set used to investigate hydrogen deuterium exchange in small molecule HRMS (Ruttkies et al. submitted). **HDXEXCH** list also contains observed deuterated species. |
| **Neurotoxicants Collection from Public Resources** | NEUROTOXINS **XLSX, CSV** (14/02/2019) CompTox **NEUROTOXINS List** | NEUROTOXINS **InChIKeys** (14/02/2019) | A list of neurotoxicants compiled from public resources, details on CompTox and Schymanski et al. (submitted). |
| **Statins Collection from Public Resources** | STATINS **XLSX, CSV** (14/02/2019) CompTox **STATINS List** | STATINS **InChIKeys** (14/02/2019) | A list of statins (lipid-lowering medications) compiled from public resources, details on CompTox. |
| **Synthetic Cannabinoids and Psychoactive Compounds** | SYNTHCANNAB **XLSX, CSV** (14/02/2019) CompTox **SYNTHCANNAB List** | SYNTHCANNAB **InChIKeys** (14/02/2019) | A list of synthetic cannabinoids and psychoactive compounds assembled from public resources, from CompTox. |

# Integration to MetFrag in place

# MassBank mapping to Dashboard Based on Web Service lookup



## MassBank Record: EA028808

Atrazine; LC-ESI-ITFT; MS2; CE: 15%; R=15000; [M+H]+

Mass Spectrum

Chemical Structure

Options
• Labels

CH$NAME: Atrazine
CH$NAME: 6-chloro-N-ethyl-N'-isopropyl-1,3,5-triazine-2,4-diamine
CH$NAME: 6-chloranyl-N4-ethyl-N2-propan-2-yl-1,3,5-triazine-2,4-diamine
CH$COMPOUND_CLASS: N/A; Environmental Standard
CH$FORMULA: C8H14ClN5
CH$EXACT_MASS: 215.0932
CH$SMILES: c1(nc(nc(n1)Cl)NCC)NC(C)C
CH$IUPAC: InChI=1S/C8H14ClN5/c1-4-10-7-12-6(9)13-8(14-7)11-5(2)3/h5H,4H2,1-3H3,(H2,10,11,12,13,14)
CH$LINK: CAS 1912-24-9
CH$LINK: CHEBI 15930
CH$LINK: KEGG C06551
CH$LINK: PUBCHEM CID:2256
CH$LINK: INCHIKEY MXWJVTOOROXGIU-UHFFFAOYSA-N
CH$LINK: CHEMSPIDER 2169
CH$LINK: COMPTOX DTXSID9020112

65

# Conclusion

- Dashboard access to data for ~875,000 chemicals

- MS-Ready data facilitates structure identification

- Related metadata facilitates candidate ranking

- Relationship mappings and chemical lists of great utility

- Dashboard and contents are one part of the solution

- New developments in progress, especially API development, will be very enabling…

# Acknowledgements

- IT Development team – especially Jeff Edwards and Jeremy Dunne
- Chris Grulke for the ChemReg system
- NERL colleagues – Jon Sobus, Elin Ulrich, Mark Strynar, Seth Newton, Alex Chao
- Emma Schymanski, LCSB, Luxembourg
- NORMAN Network and all contributors

# Contact

## Antony Williams

US EPA Office of Research and Development

National Center for Computational Toxicology

**EMAIL:** Williams.Antony@epa.gov

**ORCID**: https://orcid.org/0000-0002-2668-4821