

Delivering Transparency in Research Data: Web-based Dashboards at the National Center for Computational Toxicology at US-EPA

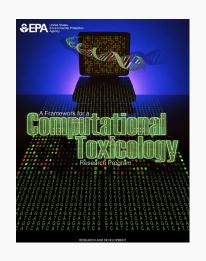
Antony Williams

U.S. Environmental Protection Agency, RTP, NC

This work was reviewed by the U.S. EPA and approved for presentation but does not necessarily reflect official Agency policy.

National Center for Computational Toxicology







- National Center for Computational Toxicology established in 2005 to integrate:
 - High-throughput and high-content technologies
 - Modern molecular biology
 - Data mining and statistical modeling
 - Computational biology and chemistry
- Currently staffed by ~60 employees as part of EPA's Office of Research and Development
- Home of ToxCast & ExpoCast research efforts
- Key partner in U.S. Tox21 federal consortium
- Multiple cross-division collaborations (e.g. NERL, OPP, OPPT)

The CompTox Portal https://comptox.epa.gov/





Environmental Topics

Laws & Regulations

About EPA

Search EPA.gov

.







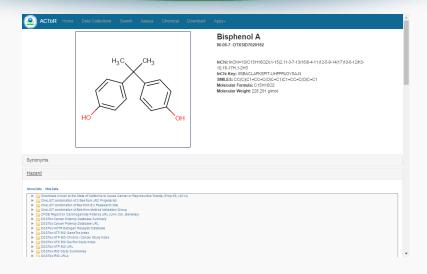


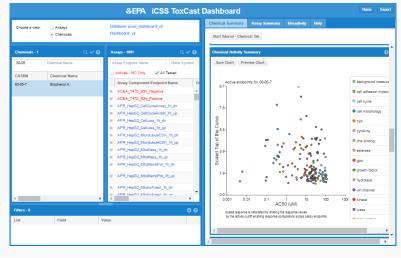


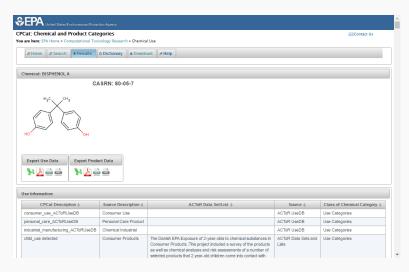


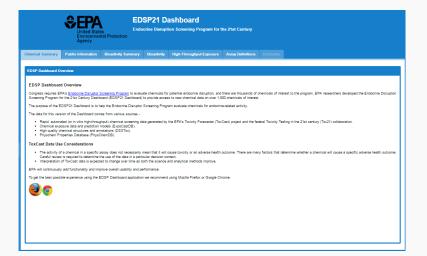
Earlier Dashboard Applications: Single architecture in development











The CompTox Portal https://comptox.epa.gov/





Environmental Topics

Laws & Regulations

About EPA

Search EPA.gov

,













The CompTox Chemicals Dashboard

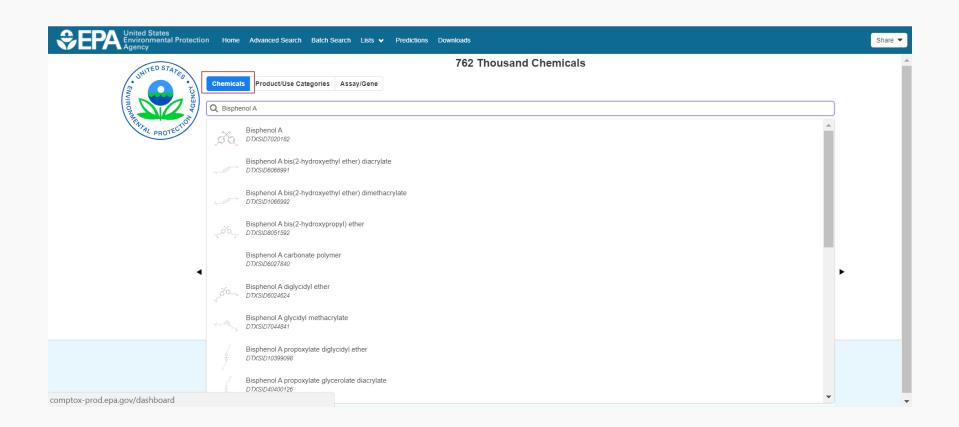


A publicly accessible website delivering access:

- ~765,000 chemicals with related property data
- Experimental and predicted physicochemical property data
- Integration to "biological assay data" for 1000s of chemicals
- Information regarding consumer products containing chemicals
- Links to other agency websites and public data resources
- "Literature" searches for chemicals using public resources
- "Batch searching" for thousands of chemicals
- DOWNLOADABLE Open Data for reuse and repurposing

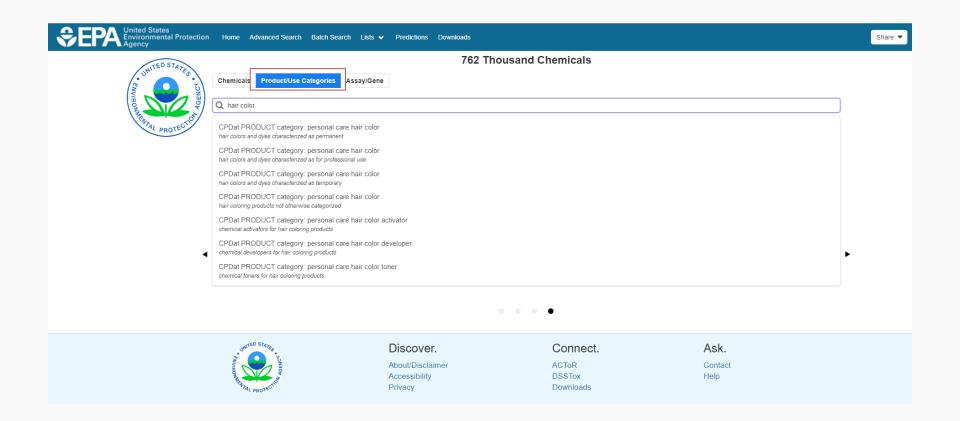
CompTox Chemicals Dashboard Chemicals





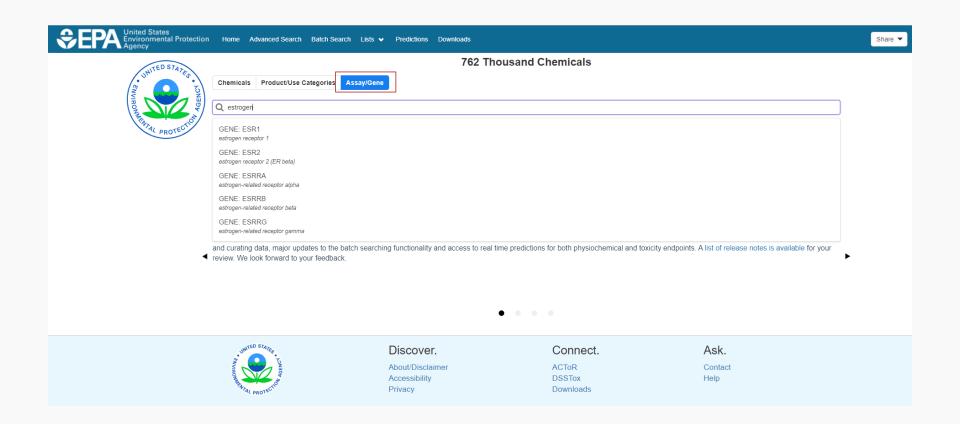
CompTox Chemicals Dashboard Products and Use Categories





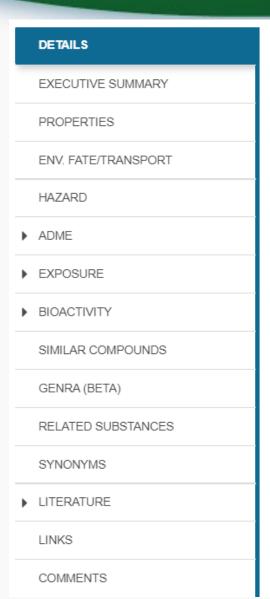
CompTox Chemicals Dashboard Assays and Genes

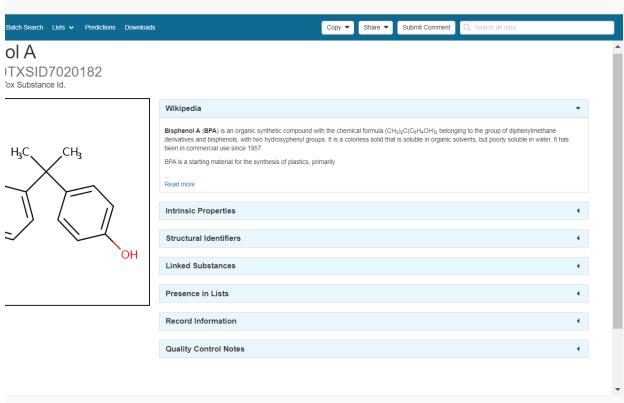




Detailed Chemical Pages

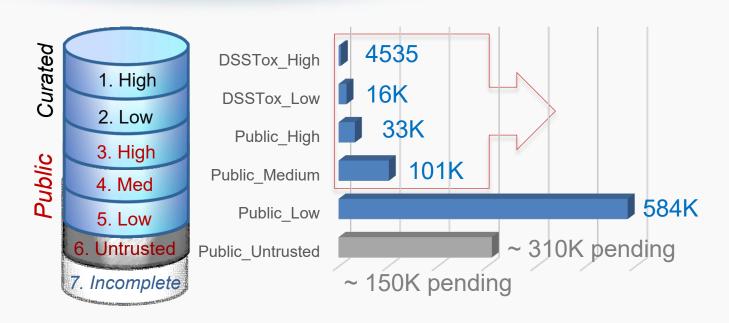






Almost 20 Years of Data... Growing with daily curation





QC Levels

DSSTox_High: Hand curated and validated

DSSTox_Low: Hand curated and confirmed using multiple public sources

Public_High: Extracted from EPA SRS and confirmed to have no conflicts in ChemID and PubChem

Public_Medium: Extracted from ChemID and confirmed to have no conflicts in PubChem

Public_Low: Extracted from ACToR or PubChem

Public_Untrusted: Postulated, but found to have conflicts in public sources

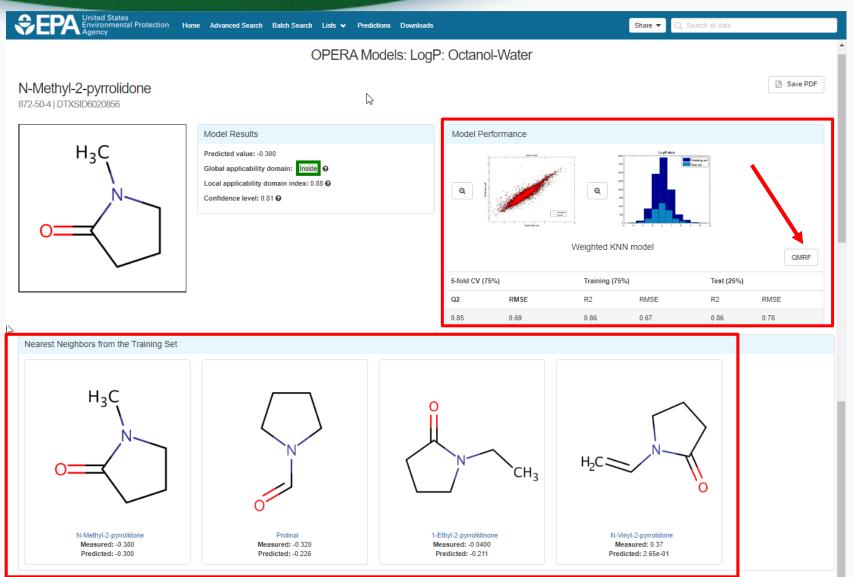
Physicochemical and Environmental Fate/Transport Properties



- Solubility
- Melting Point
- Boiling Point
- LogP (Octanol-water partition coefficient)
- Atmospheric Hydroxylation Rate
- LogBCF (Bioconcentration Factor)
- Biodegradation Half-life
- Henry's Law Constant
- Fish Biotransformation Half-life
- LogKOA (Octanol/Air Partition Coefficient)
- LogKOC (Soil Adsorption Coefficient)
- Vapor Pressure
- And more...

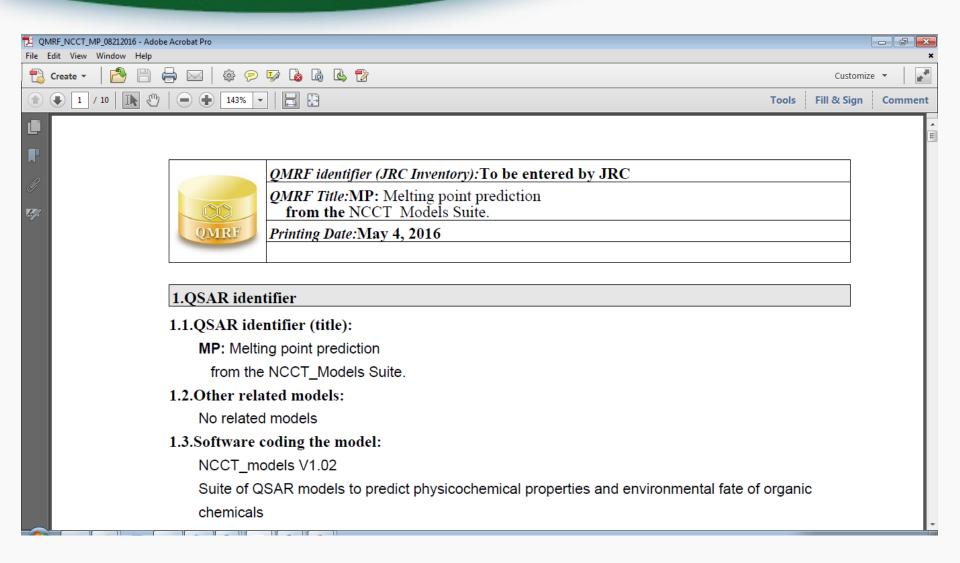
Full transparency for predictions





Transparency: QMRF Report





Developing "NCCT Models"



- Our approach to modeling:
 - Obtain high quality training sets
 - Apply appropriate modeling approaches
 - Validate performance of models
 - Define the applicability domain and model limitations
 - Use models to predict properties across our full datasets
 - Release as Open Data and Open Models

Available online



EPI Suite Data - ISIS/Base & SDF

The downloaded files are provided in "zip" format ... the downloaded file must be "un-zipped" with common utility programs such as WinZip.

... Updated September 15, 2010

Basic Instructions:

- (1) Download the zip file
- (2) Un-Zip the file

<u>NOTE</u> ... zipped files extract to Folders containing the individual data files ... Folders named EPI_ISIS_Data and EPI_SDF_Data

<u>Substructure Searching Files:</u>

ISISTM/Base & SD Files of the EPI Suite Program Experimental Data Files are now available ... The ISISTM/Base files require the commercial program for use ... The SD Files can be imported into other commercial chemical structure programs (such as ChemFinder).

... Click here to download EPI_ISIS_Data.zip ... (about 11 MB)

... Click here to download EPI_SDF_Data.zip ... (about 10 MB)

NOTE ... EPI Suite Data Files (some in Excel, Text, Word format) available at:

http://esc.syrres.com/interkow/EpiSuiteData.htm

We Curated These Public Data to Build Prediction Models



Public data should be curated prior to modeling **Different Compounds**

Mol Block	S CAS	S NAME	Smiles
H ₃ C OH	000076-43-7	FLUOXYMESTERONE	H 3 C
H ₃ C CH ₃	000077-99-6	1,1,1-TRIS(HYDROXYMETHYL)PROPANE	HO OH CH3
он он	000079-60-7	CORTISONE-9A-FLUORO	O C C C C C C C C C C C C C C C C C C C
NH ₂	000082-38-2	DISPERSE RED 9	H3 C - N

Duplicate Structures



Structure	Formula <	FW <	CAS <	NAME <	MP <	EstMP <	ErrorMP <
OH OH OH	с ₃ н ₆ о ₃	90.0779	000050-21-5	LACTIC ACID	1.6800000000000 00e+001	2.2660000000000 00e+001	5.860000000000 00e+000
OH OH OH	с ₃ н ₆ о ₃	90.0779	000079-33-4	L-LACTIC ACID	5.300000000000 00e+001	2.2660000000000 00e+001	-3.03400000000 000e+001
O OH OH	с ₃ н ₆ о ₃	90.0779	000598-82-3	A-HYDROXYPROPIONIC ACID	1.8000000000000 00e+001	2.2660000000000 00e+001	4.6600000000000 00e+000
O OH OH	с ₃ н ₆ о ₃	90.0779	010326-41-7	D-LACTIC ACID	5.2800000000000 00e+001	2.2660000000000 00e+001	-3.01400000000 000e+001

Covalent Halogens



Mol Block	S CAS	S NAME	Smiles
H ₃ C CH ₃	000056-93-9	BENZYL TRIMETHYL AMMONIUM CHLORIDE	CH ₃ H ₃ C - N CH ₃ CI
H ₃ C CH ₃ CH ₃ CH ₃	000068-05-3	TETRAETHYL AMMONIUM IODIDE	H ₃ C CH ₃ CH ₃ CH ₃
H ₃ C CH ₃	000071-91-0	TETRAETHYL AMMONIUM BROMIDE	H ₃ C CH ₃ Br CH ₃

Curation to QSAR Ready Files



Property	Initial file	Curated Data	Curated QSAR ready
AOP	818	818	745
BCF	685	618	608
BioHC	175	151	150
Biowin	1265	1196	1171
ВР	5890	5591	5436
HL	1829	1758	1711
KM	631	548	541
KOA	308	277	270
LogP	15809	14544	14041
MP	10051	9120	8656
PC	788	750	735
VP	3037	2840	2716
WF	5764	5076	4836
WS	2348	2046	2010

LogP dataset:15,809 structures



- CAS Checksum: 12163 valid, 3646 invalid (>23%)
- Invalid names: 555
- Invalid SMILES 133
- Valence errors: 322 Molfile, 3782 SMILES (>24%)
- Duplicates check:
 - 31 DUPLICATE MOLFILES
 - 626 DUPLICATE SMILES
 - 531 DUPLICATE NAMES
- SMILES vs. Molfiles (structure check)
 - 1279 differ in stereochemistry (~8%)
 - 362 "Covalent Halogens"
 - 191 differ as tautomers
 - 436 are different compounds (~3%)

OPERA Predicted Properties



An automated curation procedure for addressing chemical errors and inconsistencies in public datasets used in QSAR modelling

K. Mansouri, C. M. Grulke, A. M. Richard, R. S. Judson & A. J. Williams



Journal of Cheminformatics

RESEARCH ARTICLE

Open Access

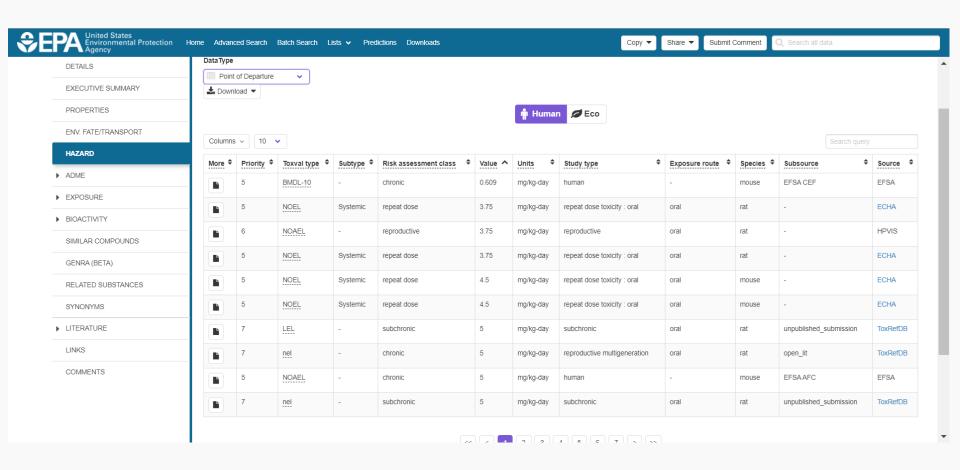
OPERA models for predicting physicochemical properties and environmental fate endpoints

Kamel Mansouri 1,2,3* , Chris M. Grulke 1, Richard S. Judson 1 and Antony J. Williams 1

OPERA Models: https://github.com/kmansouri/OPERA

Other Data: Human and Ecological Chemical Hazard Data





Hazard Data from "ToxVal_DB"



- ToxVal Database contains following data:
 - -30,050 chemicals
 - -772,721 toxicity values
 - -29 sources of data
 - -21,507 sub-sources
 - -4585 journals cited
 - -69,833 literature citations

If only the data were easy to extract...

Some ways to do it poorly...



Name	RT	m/z	Formula	m/z Diff (ppm)	Mass score	SMILES
Parent	3.13	455.2926	C27H38N2O4	-3.48		N#CC(CCCN(C)CCc1ccc(OC)c(OC)c1)(C(C)C)c2ccc(OC)c(OC)c2
M6 -164	2.26	291.2077	C17H26N2O2	− 1.37	429	N(C)CCCC(C#N)(C(C)C)cIccc(OC)c(OC)cI
MI6 -I4	3.06	441.2743	C26H36N2O4	2.41	534	c1cc(CCNCCCC(C#N)(C(C)C)c2ccc(OC)c(OC)c2)cc(OC)c1OC
MI4 +I6	2.92	471.2866	C27H38N2O5	−1.59	476	N#CC(CCCN(C)CC(O)clccc(OC)c(OC)cl)(C(C)C)c2ccc(OC)c(OC)c2
M9 -14	2.78	441.2761	C26H36N2O4	−1.7	590	C(#N)C(CCCN(C)CCc1ccc(OC)c(OC)c1)(C(C)C)c2ccc(O)c(OC)c2
MII -14	2.84	441.2742	C26H36N2O4	2.57	473	Oclccc(CCN(C)CCCC(C#N)(C(C)C)c2ccc(OC)c(OC)c2)ccIOC
M12 +2	2.87	457.2707	C26H36N2O5	-0.93	570	O(C)clcc(ccclOC)C(C#N)(CCCNCC(O)c2ccc(OC)c(OC)c2)C(C)C
M5 -178	2.2	277.1894	C16H24N2O2	7.84	419	C(C)(C)C(C#N)(CCCN)c1ccc(OC)c(OC)c1
M8 +2	2.67	457.2708	C26H36N2O5	-I.18	581	OC(CN(C)CCCC(C#N)(C(C)C)clccc(OC)cl)c2ccc(O)c(OC)c2
MI5 -I4	2.92	441.2743	C26H36N2O4	2.23	614	Oclccc(CCN(C)CCCC(C#N)(C(C)C)c2ccc(OC)c(OC)c2)cclOC
M2 -259	0.73	196.1326	CI IH I7NO2	5.91	618	cI(CCNC)ccc(OC)c(cI)OC
MI0 -28	2.8	427.2617	C25H34N2O4	−4.79	487	c1(OC)cc(ccc1OC)C(C#N)(CCCNCCc2ccc(O)c(OC)c2)C(C)C
M7 +2	2.46	457.2717	C26H36N2O5	-3.23	492	c1cc(CCN(O)CCCC(C#N)(C(C)C)c2ccc(OC)c(OC)c2)cc(OC)c1OC
MI7 +16	3.21	471.2853	C27H38N2O5	1.19	534	N#CC(CCCN(C)CCc1ccc(OC)c(OC)c1)(c2ccc(OC)c(OC)c2)C(C)(C)OC
M4 -178	1.86	277.1927	C16H24N2O2	-3.8	444	COcIcc(cccIO)C(C#N)(CCCNC)C(C)C
MI -289	0.44	166.0858	C9H1 INO2	5.93	136	c1(CCNC)ccc(=O)c(c1)=O
MI3 -16	2.88	439.2603	C26H34N2O4	-I.43	367	c1cc(CC=NCCCC(C#N)(C(C)C)c2ccc(OC)c(OC)c2)cc(OC)c1OC

How do I extract structures?



Copy-Paste doesn't work

```
cI(CCNC)ccc(=O)c(cI)=O

cIcc(CC=NCCCC(C#N)(C(C)C)c2ccc(OC)c(OC)c2)cc(OC)cIOC

c1(CCNC)ccc(*O)c(c1)*0
c1cc(CC*NCCCC(C#N)(C(C)C)c2ccc(OC)c(OC)c2)cc(OC)c1OC
```

• This is not the way a publisher should deliver chemistry. But this is on the AUTHOR! 25

Trust automatic extraction of Structure Drawings???

260.1637

260.1651

5.33

Table 2. Selection of fragments that help in the M16-16 metabolite structure elucidation



Sub. obs. m/z	Sub. cal. m/z	Sub. m/z diff. ppm	Substrate	Metabolite	Δ	Met. obs. m/z	Met. calc. m/z	Met. m/z diff. ppm
150.0664	150.0681	11.42		100	+0	150.0670	150.0681	7.25
165.0869	165.0916	28.22			+0	165.0892	165.0916	14.30

+0

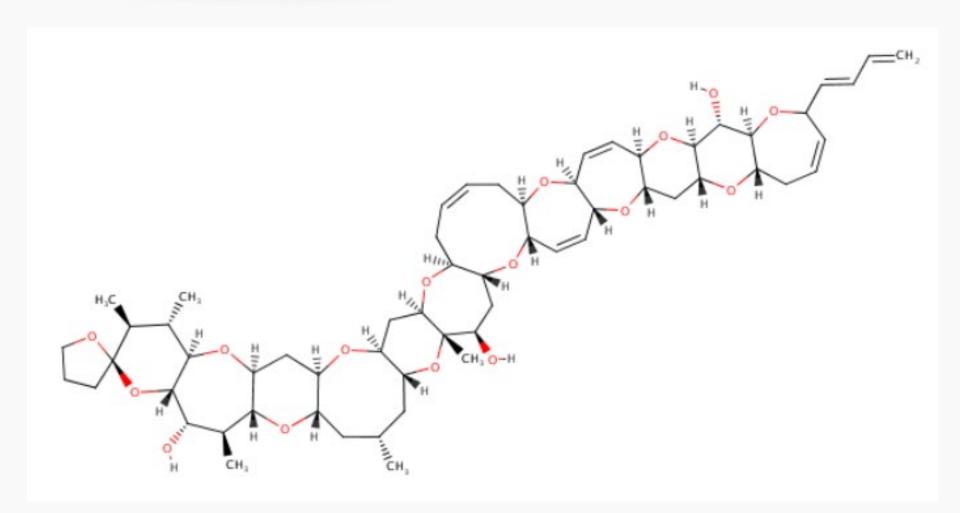
260.1652

260.1651

-0.50

Try hand-drawing Algal Toxins!

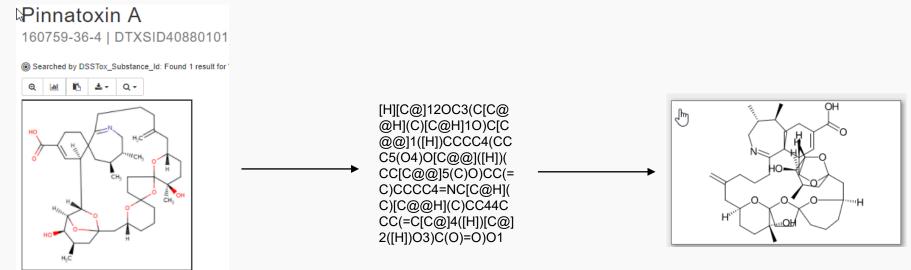




Think of files in multiple formats!



 SMILES are hyper-dependent on good layout algorithms. It's not easy!



 We publish our papers with Excel files and SDF files as supp info, on FigShare, and on the dashboard as data collections

Not just chemical "structures"



 Chemicals in commerce, of interest to the EPA, are not all easily represented by structures

- Different chemical substances supported
 - "UVCB chemicals" Unknown or Variable Composition,
 Complex Reaction Products and Biological Materials
 - Homologous series as Markush Structures

Example PFAS-UVCBs



0 related chemical structures with this substance

Ethene, tetrafluoro-, oxidized, polymd., ...
DTXSID: DTXSID00108075
CASRN: 274917-96-3

0 related chemical structures with this substance

Sulfonamides, C4-8-alkane, perfluoro, ...
DTXSID: DTXSID00108095
CASRN: 160901-25-7

0 related chemical structures with this substance

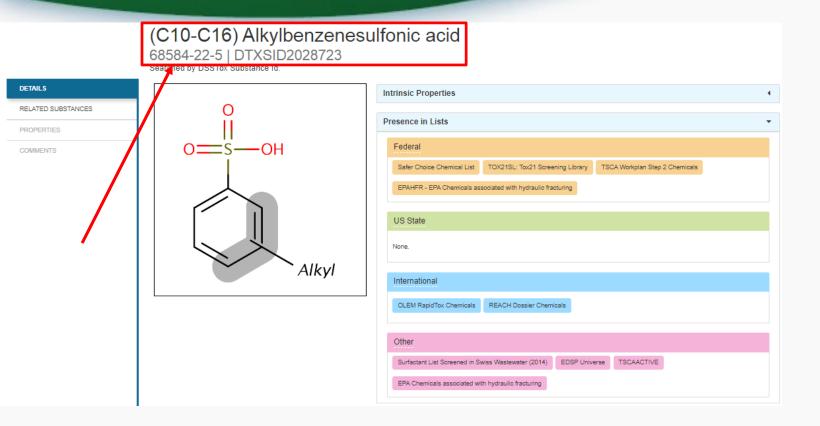
1-Propene, 1,1,2,3,3,3-hexafluoro-, pol... DTXSID: DTXSID00108732 CASRN: 149935-01-3

Ethene, tetrafluoro-, oxidized, polymd., reduced, decarboxylated, C6 fraction 274917-96-3 | DTXSID00108075

1-Propene, 1,1,2,3,3,3-hexafluoro-, polymer with 1,1-difluoroethene, ethene, 1,1,2,2-tetrafluoroethene and 1,1,2-trifluoro-2-(trifluoromethoxy)ethene

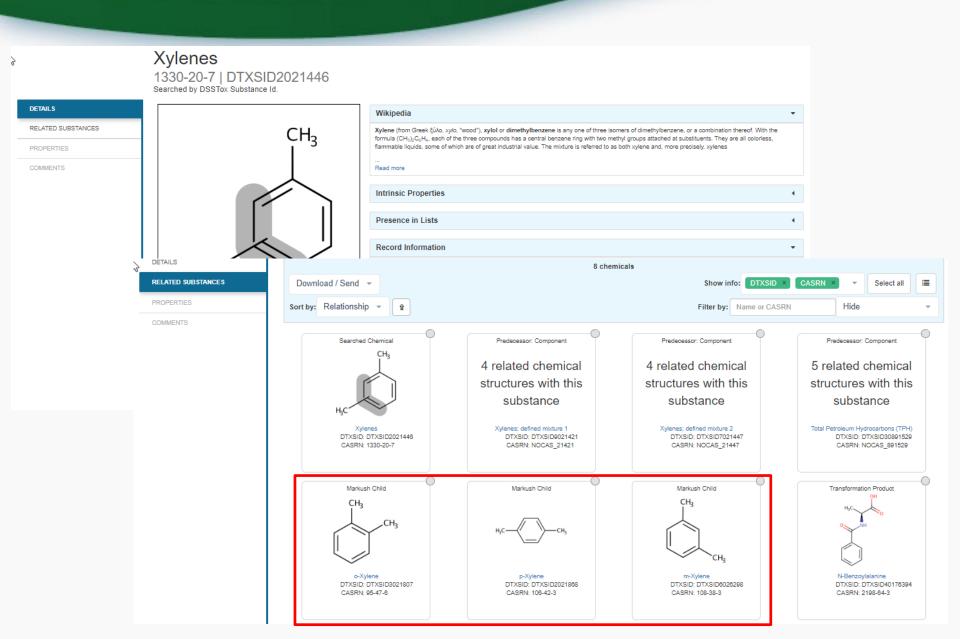
UVCB Chemicals





Markush Structures





Environmental Chemistry: More about Names and CASRNs



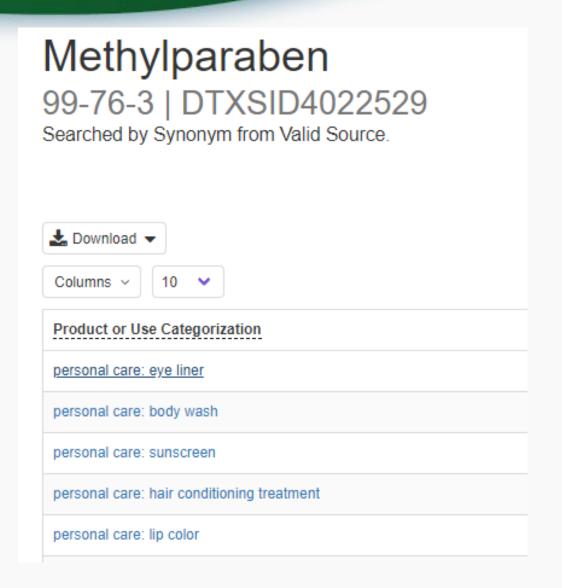
 For EPA most chemicals are reported as text – chemical names and CAS Numbers

Attachment D (Method 3) SIM quantitation ions and qualifiers for internal standards, references method analysis, and surrogates

Name of Compound	CAS No.	Ouantitation Ion	Oualifier Ions
Phenol-d6 (SS)	13187-88-3	99	71, 42
Phenol	108-95-2	94	66
1,4-Dichlorobenzene	106-46-0	146	111. 75, 50
Acetophenone	98-86-2	105	77, 51, 120
Acenaphthene-d10 (IS)	15067-26-2	162	160, 80
p-Cresol	106-44-5	107	108, 77
Isophorone	78-59-1	82	138, 54
Camphor	76-22-2	95	81, 108, 152
Isoborneol	124-76-5	95	110, 121, 136
Menthol	89, 78, 1	71	81, 123, 138
Naphthalene	91-20-3	128	102, 51
Methyl salicilate	119-36-8	120	92, 152, 65 ₃₃

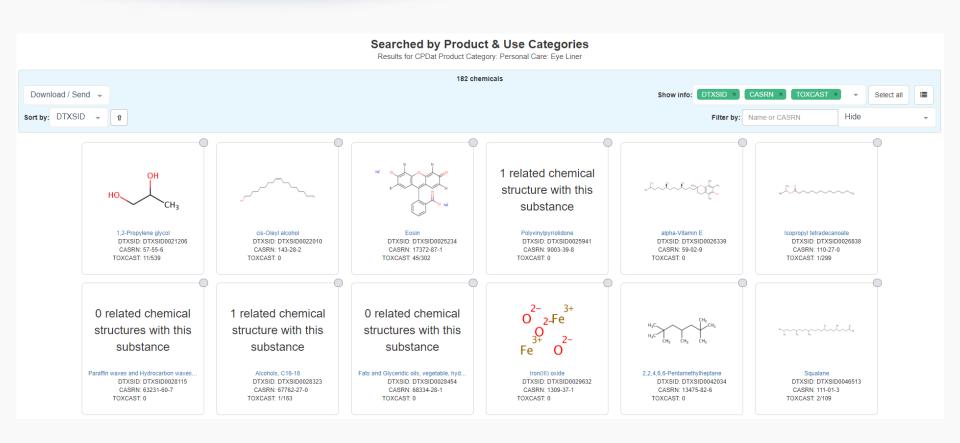
Chemicals in product SDS sheets are commonly UVCBs





182 chemicals in Personal Care: Eye Liner Category





CASRNs can be problematic...



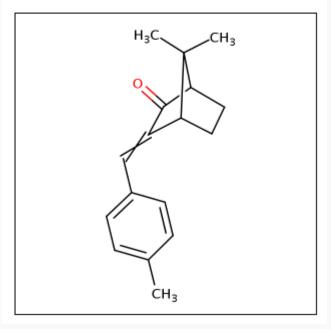
Name of Compound	CAS No.
Phenol-d6 (SS)	13187-88-3
Phenol	108-95-2
1,4-Dichlorobenzene	106-46-0
Acetophenone	98-86-2
Acenaphthene-d10 (IS)	15067-26-2
p-Cresol	106-44-5
Isophorone	78-59-1
Camphor	76-22-2
Isoborneol	124-76-5
Menthol	89, 78, 1
Naphthalene	91-20-3
Methyl salicilate	119-36-8

Active vs Deleted CASRN Also "Alternates"



Enzacamene

36861-47-9 | DTXSID8047896 Searched by Approved Name.



Synonym	Quality
Enzacamene	Valid
7,7-Dimethyl-3-[(4-methylphenyl)methylidene]bicyclo[2.2.1]heptan-2-one	Valid
Bicyclo[2.2.1]heptan-2-one, 7,7-dimethyl-3-[(4-methylphenyl)methylene]-	Valid
36861-47-9 Active CA 8-FN	Valid
Bicyclo[2.2.1]heptan-2-one, 1,7,7-trimethyl-3-[(4-methylphenyl)methylene]-	Valid
EINECS 253-242-8	Other
Eusolex 6300	Other
Uvinul MBC 95	Other
Parsol 5000	Other

ONII-8I3XWY40L9	Other
4-Methylbenzylidenecamphor	Other
p-Methylbenzylidenecamphor	Other
38102-82-4 Deleted GAS-RIN	Deleted
84055-65-2 Detected CA 8-RN	Deleted

Tricky mapping by CASRN This one has 316 Deleted CASRN



CAS Registry Number: 25068-38-6

(C₁₅ H₁₆ O₂ . C₃ H₅ CI O)_X

Phenol, 4,4'-(1-methylethylidene)bis-, polymer with 2-(chloromethyl)oxirane

Polymer

Polymer Class Terms: Epoxy resin

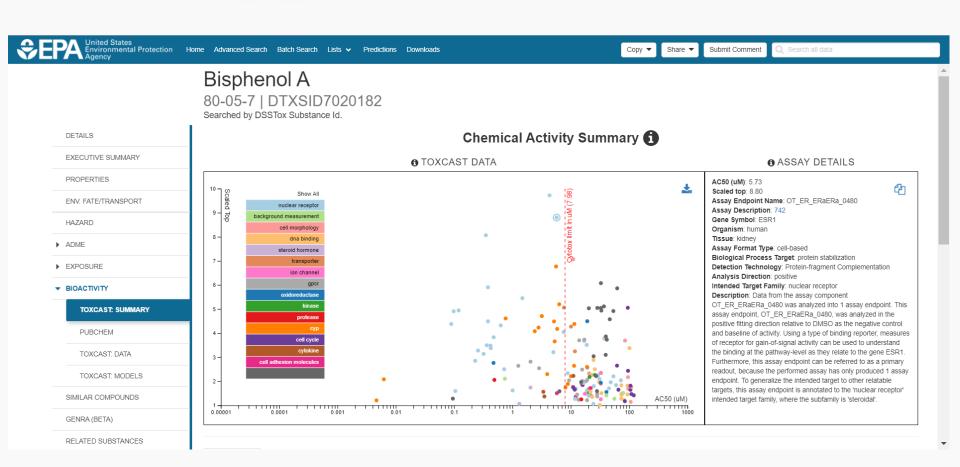
Alternate CAS Registry Numbers: 26402-79-9

Deleted CAS Registry Numbers: 1336-88-5, 1337-15-1, 8000-31-5, 9015-99-0, 9049-54-1, 9050-21-9, 9081-91-8, 9083-76-5, 9084-94-0, 9086-62-8, 9087-26-7, 9087-76-7, 11097-80-6, 11098-13-8, 11098-40-1, 11100-23-5, 11108-41-1, 11120-31-3, 11121-19-0, 11126-36-6, 20232-24-0, 35038-60-9, 7, 37270-82-9, 37291-75-1, 37293-07-5, 37294-18-1, 37305-82-1, 37307-45-2, 37317-45-6, 37325-21-6, 37338-63-9, 37342-17-9, 37345-34-9, 37348-56-4, 37348-57-5, 37357-73-6, 37360-93-3, 39277-59-3, 39288-99-8, 39296-08-7, 39296-09-8, 39296-11-2, 39296-15-6, 39315-77-0, 39349-91-2, 39354-86-4, 39362-25-9, 39362-45-3, 39373-81-4, 39378-29-5, 39378-55-7, 39389-49-6, 39405-18-0, 39412-57-2, 39419-66-4, 39453-22-0, 39454-54-1, 39454-69-8, 39470-62-7, 42612-34-0, 42618-03-1, 50642-36-9, 50642-55-2, 50642-78-9, 51158-20-4, 51273-81-5, 51329-73-8, 51393-99-8, 51394-03-7. 51553-00-5. 52011-87-7. 52038-45-6. 52051-70-4. 52051-82-8. 52052-16-1. 52232-05-0. 52232-75-4. 52276-55-8. 52365-33-0. 52519-66-1. 52519-67-2. 52627-94-8. 52907-38-7. 53027-88-6. 53127-14-3. 53200-30-9. 53238-86-1. 53238-87-2. 53239-67-1. 53239-68-2. 53570-97-1. 53570-98-2. 53681-78-0. 53858-93-8. 54018-73-4. 54352-05-5. 55464-96-5. 55584-55-9. 55585-07-4. 55818-73-0. 56258-35-6. 56449-43-5. 56509-48-9. 57107-66-1, 57284-90-9, 57534-21-1, 57693-04-6, 58052-05-4, 58128-38-4, 58392-89-5, 58392-92-0, 58516-14-6, 58572-71-7, 59029-19-5, 59459-14-2, 59473-30-2, 59948-36-6, 60202-19-9, 60267-31-4, 60382-89-0, 60606-56-6, 60800-54-6, 60831-77-8, 60894-16-8, 61036-82-6, 61287-42-1, 61356-27-2. 61711-38-4. 61763-30-2. 61991-18-2. 62169-28-2. 62169-29-3. 62601-75-6. 62601-76-7. 62887-23-4. 63055-40-3. 63172-55-4. 63799-24-6. 63993-57-7. 63993-58-8. 64086-14-2. 64086-16-4. 64176-52-9. 64176-61-0. 64176-66-5. 64177-03-3. 65233-49-0. 65931-38-6. 65931-39-7. 66995-96-8. 67185-62-0. 68821-97-6. 69899-40-7. 70179-83-8. 70213-44-4. 70726-45-3. 71965-91-8. 72514-40-0. 73413-19-1. 74504-20-4. 74564-76-4. 75831-44-6, 78564-77-9, 79585-43-6, 80702-61-0, 81458-12-0, 81843-57-4, 81843-58-5, 81855-87-0, 82197-12-4, 82197-46-4, 83202-85-1, 84286-97-5, 84683-04-5, 84931-29-3, 85537-69-5, 86090-60-0, 88385-37-9, 88528-19-2, 88651-18-7, 89750-00-5, 91727-28-5, 91727-29-6, 92481-37-3, 95327-25-6, 96420-31-4, 96510-68-8, 97568-16-6, 97709-01-8, 99400-50-7, 101027-12-7, 102256-87-1, 103599-13-9, 103599-14-0, 104364-97-8, 104491-99-8, 105521-57-1, 106207-08-3, 106856-89-7, 107991-47-9, 108556-05-4, 108728-21-8, 110158-22-0, 111367-08-9, 111517-59-0, 114013-37-5, 115902-32-4, 117216-90-7, 117313-45-8, 117786-92-2, 118340-04-8, 120146-74-9, 120797-43-5, 121181-85-9, 121273-37-8, 121547-73-7, 123939-44-6, 125147-87-7, 127176-80-1, 127176-81-2, 128281-71-0, 132822-20-9, 132893-73-3, 135976-90-8, 137545-29-0, 138157-20-7, 138361-18-9, 139554-29-3. 142540-11-2. 144046-24-2. 144046-25-3, 144855-66-3, 149013-58-1, 150825-32-4, 157321-42-1, 157481-46-4, 158725-45-2, 160674-45-3, 161937-12-8. 162031-55-2. 167972-06-7. 168042-08-8. 179607-24-0. 183581-68-2. 183890-12-2. 187619-11-0. 188448-56-8. 189282-49-3. 191606-83-4, 220090-06-2, 222835-65-6, 222835-66-7, 222835-68-9, 222835-69-0, 222835-70-3, 222835-72-5, 222835-74-7, 222835-77-0, 309945-96-8, 339530-81-3, 353239-57-3, 367523-08-8, 383889-26-7, 383889-27-8, 395069-05-3, 470462-49-8, 681001-41-2, 848887-61-6, 913745-83-2, 917483-69-3, 922728-11-8, 934588-09-7, 945610-97-9, 950907-45-6, 1033821-54-3, 1034342-45-4, 1068160-75-7, 1082736-74-0, 1096473-97-0, 1114797-08-8. 1189565-70-5. 1190235-62-1. 1190729-68-0. 1192045-32-1. 1195324-26-5. 1196030-95-1. 1198291-96-1. 1199811-18-1. 1203835-26-0. 1206700-05-1, 1228639-00-6, 1245563-83-0, 1271727-39-9, 1300093-58-6, 1300102-07-1, 1305321-17-8, 1338071-08-1, 1446691-72-0, 1450839-98-1. 1620807-39-7. 1641551-32-7. 1807886-28-7. 1815624-46-4. 1815624-47-5

In Vitro Bioassay Screening

ToxCast and Tox21

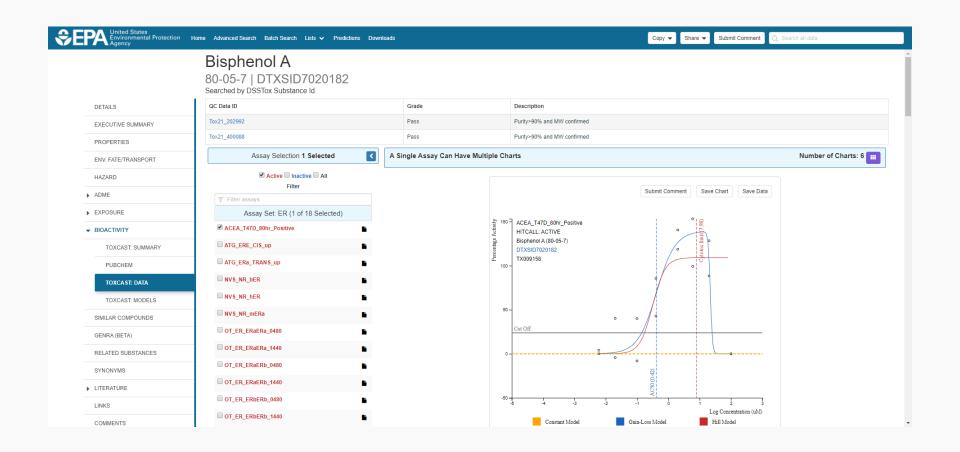




In Vitro Bioassay Screening

ToxCast and Tox21





Bioactivity: Downloadable Data

https://www.epa.gov/chemical-research/exploring-toxcast-data-downloadable-data



Exploring ToxCast Data: Downloadable Data

The results after processing through the Pipeline are available on the <u>ToxCast Dashboard</u>, and for most users EPA recommends accessing the data there.

- ToxCast Chemicals
- ToxCast Assays

ToxCast Data and Information

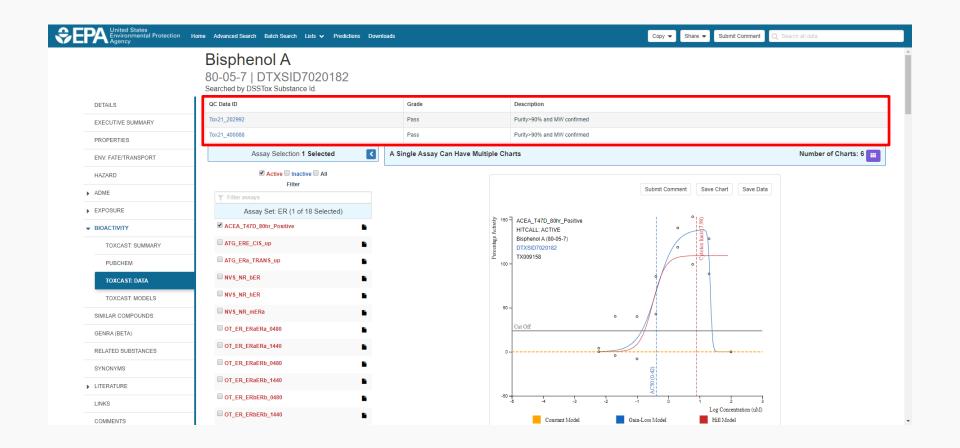
- ToxCast & Tox21 Summary Files. Data for a single chemical endpoint pair for thousands of chemicals and assay endpoints for 20 variables such as the activity or hit call, activity concentrations, whether the chemical was tested in a specific assay, etc.
 - o <u>Download ToxCast Summary Information</u>
 - Download ReadMe
- ToxCast & Tox21 Data Spreadsheet. A spreadsheet of EPA's analysis of the chemicals screened through ToxCast and the Tox21 collaboration which includes EPA's activity calls from the screening of over 1,800 chemicals.
 - Download Data
 - Download ReadMe
- ToxCast Data Pipeline R Package. The R computer programming package used to process and model all EPA ToxCast and Tox21 chemical screening data. The files include the R programming package as well as documents that provide overviews of the data analysis pipeline used and the R package. Users will need experience with R to use these files.
 - <u>Download Package</u>
 - TCPL Overview

Resources

- <u>Toxicity Forecaster (ToxCast)</u>
 <u>Fact Sheet</u>
- ToxCast Publications
- ToxCast Citation
- About ToxCast

ToxCast/Tox21 Data Analytical QC of the chemicals





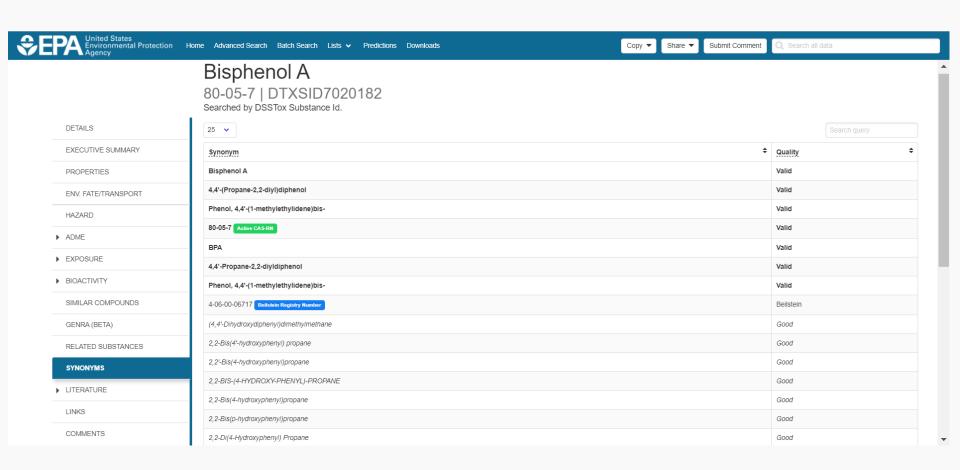
Access to Analytical QC Data





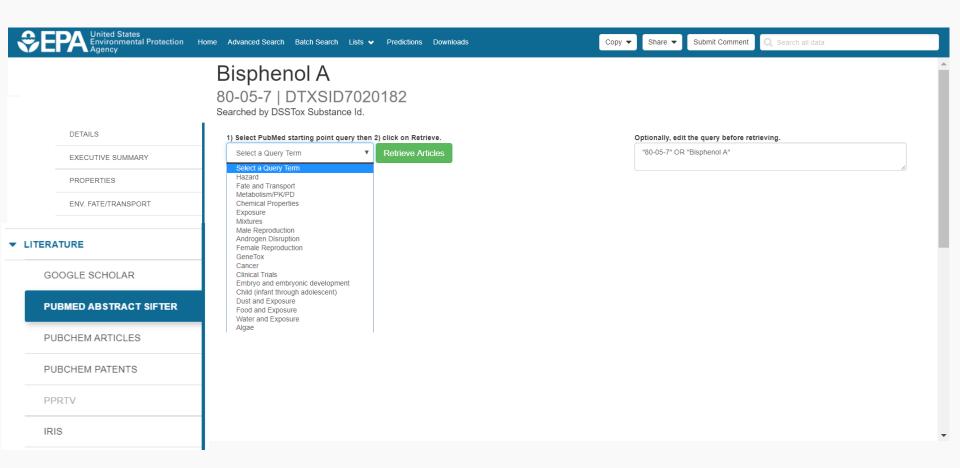
Names and CASRNs to Support Searches





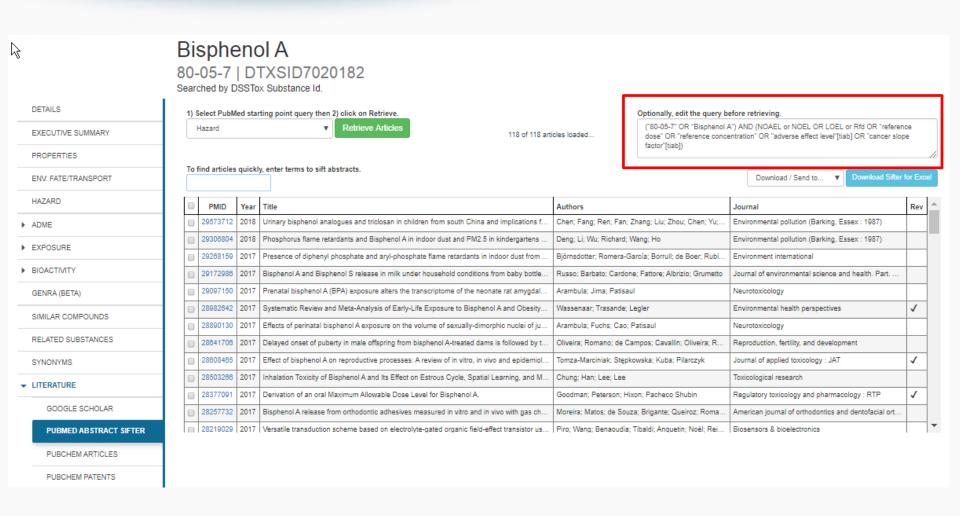
Literature Searches - Querying 28 Million PubMed abstracts





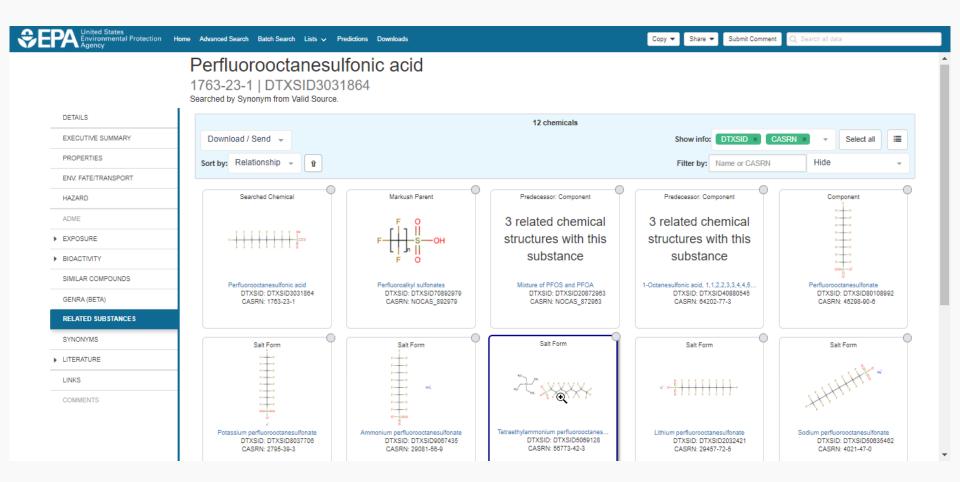
Abstract Sifter - Querying 28 Million PubMed abstracts





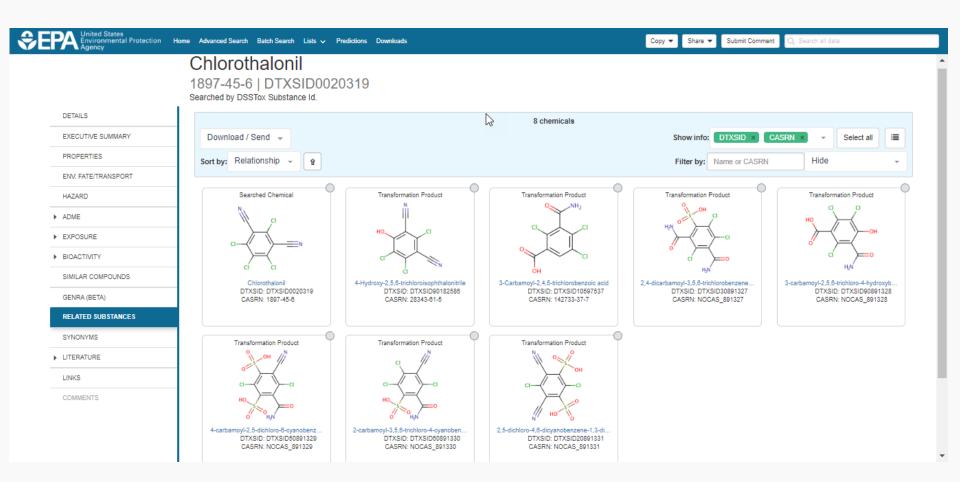
Relationships in the data





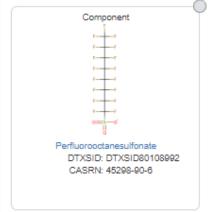
Transformation Products





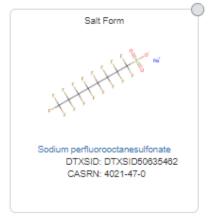
7 salt forms of PFOS (and the ion itself)

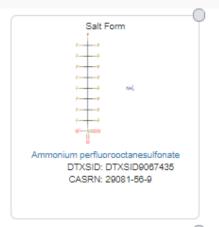


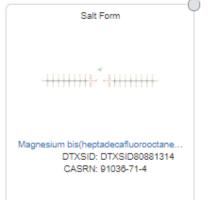


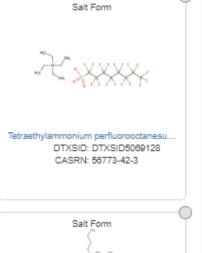


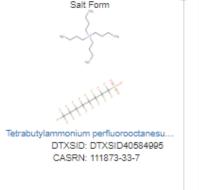








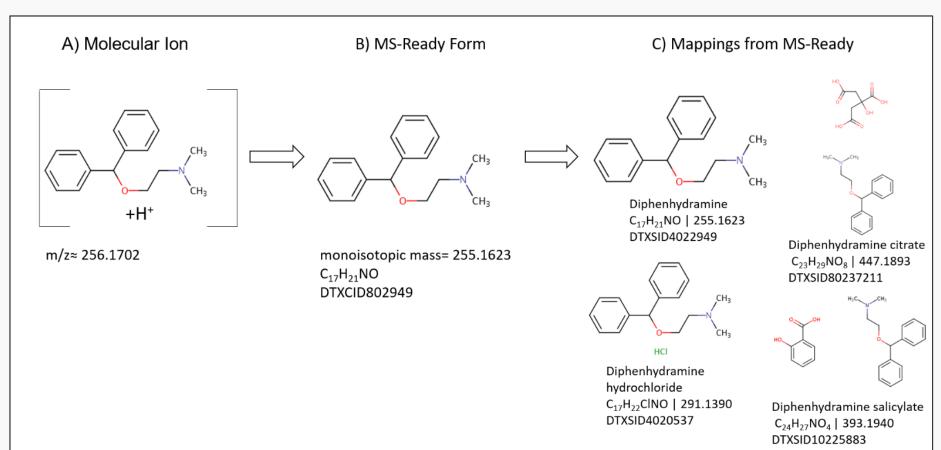




Using data relationships

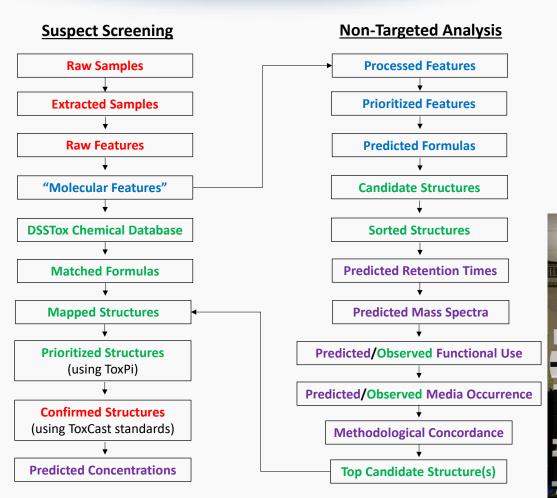


 We have purposely built relationships in the data. Specifically, "MS-Ready mappings"



"MS-Ready": Suspect Screening and Non-Targeted Analysis Workflow





Color Key

Red = Analytical Chemistry

Blue = Data Processing & Analysis

Purple = Mathematical & QSPR Modeling

Green = Informatics & Web Services



Batch Searching



 Singleton searches are useful but we need to search thousands of masses, formulae, names, InChls and CASRNs!

Typical questions

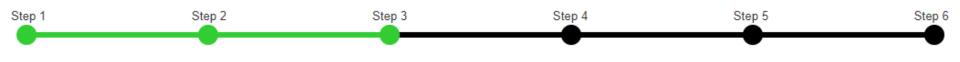
- What chemicals can I get for 5000 CAS Numbers?
- Can I get predicted properties for 1000 chemicals?
- What is the list of chemicals for the formula C_xH_yO_z?
- What is the list of chemicals for a mass +/- error ?
- Can I get chemical lists in Excel files? In SDF files?

Batch Searching

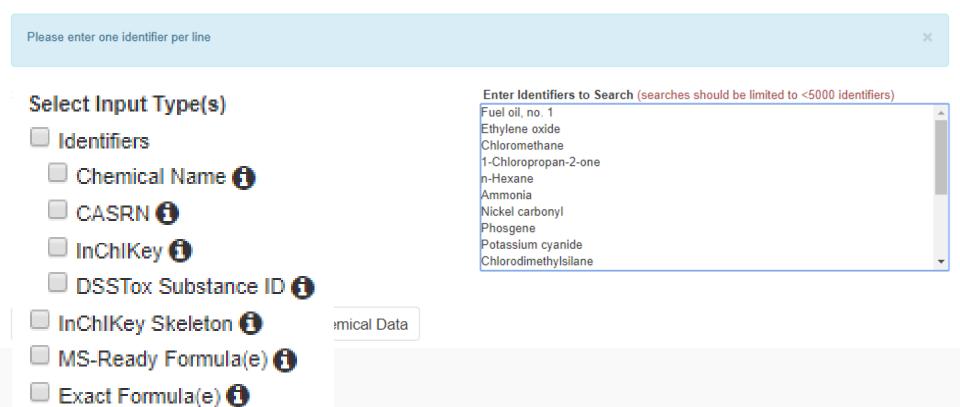
Monoisotopic Mass



Batch Search 2

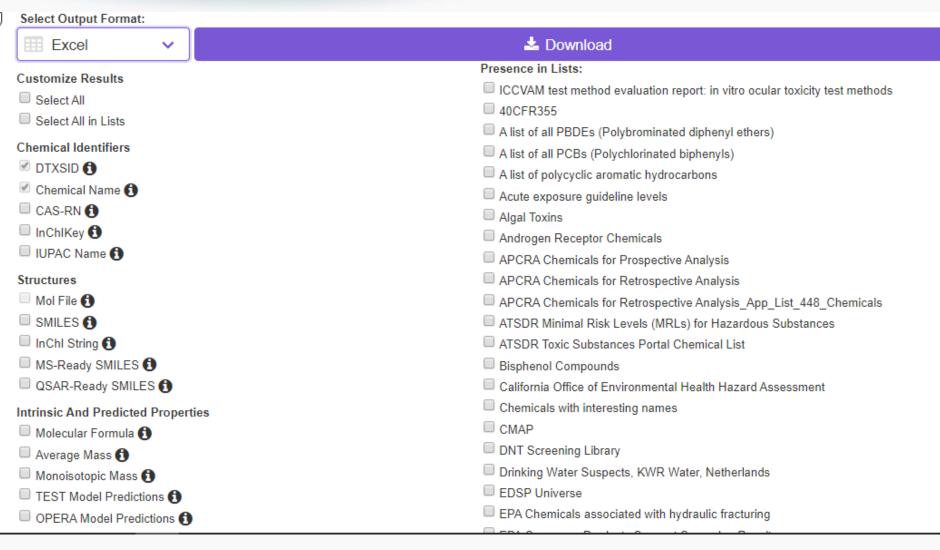


Step Three: Select Download Data or Display Chemicals



Batch Searching





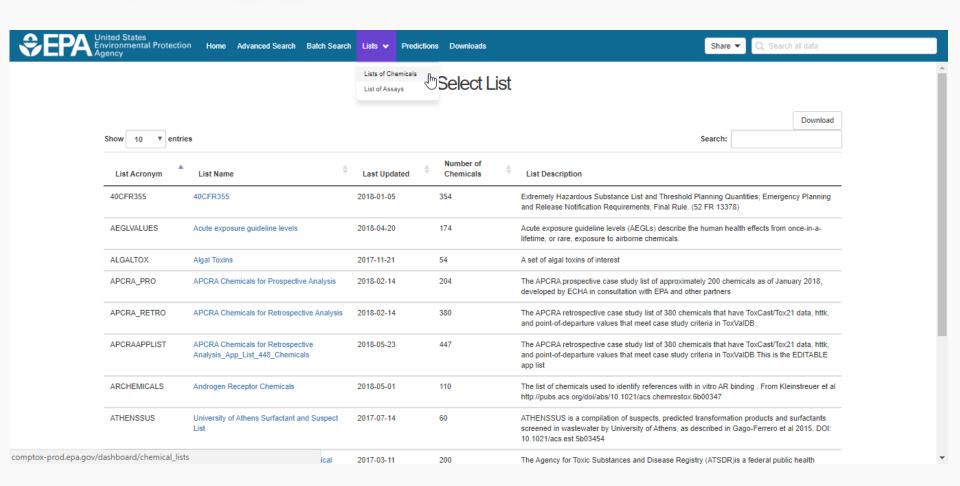
Excel Output



	Α	В	С	D	Е	F	G	Н
1	INPUT	FOUND_BY	DTXSID	PREFERRED_NAME	EXPOCAST	EXPOCAST	NHANES	TOXVAL_D
2	1445-75-6	CAS-RN	DTXSID5024051	Diisopropyl methylpho:	2.09e-08	Υ	-	Υ
3	50-00-0	CAS-RN	DTXSID7020637	Formaldehyde	1.32e-06	Υ	-	Υ
4	107-06-2	CAS-RN	DTXSID6020438	1,2-Dichloroethane	4.9e-06	Υ	-	Υ
5	57-12-5	CAS-RN	DTXSID6023991	Cyanide	-	-	-	Υ
6	7550-45-0	CAS-RN	DTXSID8042476	Titanium tetrachloride	_	-	-	Υ
7	79-01-6	CAS-RN	DTXSID0021383	Trichloroethylene	7.27e-06	Υ	-	Υ
8	121-82-4	CAS-RN	DTXSID9024142	Cyclonite	6.72e-08	Υ	-	Υ
9	108-05-4	CAS-RN	DTXSID3021431	Vinyl acetate	8.3e-05	Υ	_	Υ
10	7803-51-2	CAS-RN	DTXSID2021157	Phosphine	-	-	-	Υ
11	122-66-7	CAS-RN	DTXSID7020710	1,2-Diphenylhydrazine	1.49e-07	Υ	_	Υ
12	101-77-9	CAS-RN	DTXSID6022422	4,4'-Methylenedianiline	6.08e-06	Υ	-	Υ
13	14017-34-6	CAS-RN	DTXSID90161250	Selenium difluoride	-	-	-	-
14	75-44-5	CAS-RN	DTXSID0024260	Phosgene	-	-	-	Υ
15	621-64-7	CAS-RN	DTXSID6021032	N-Nitrosodipropylamine	4.55e-07	Υ	_	Υ
16	75-09-2	CAS-RN	DTXSID0020868		2.02e-06	Υ	-	Υ
17	100-41-4	CAS-RN	DTXSID3020596	Ethylbenzene	8.32e-05	Υ	-	Υ
18	7440-28-0	CAS-RN	DTXSID2036035	Thallium	-	-	-	Υ
19	108-88-3	CAS-RN	DTXSID7021360	Toluene	8.61e-05	Υ	-	Υ
20	111-44-4	CAS-RN	DTXSID9020168	Bis(2-chloroethyl) ethe	2.82e-07	Υ	-	Υ
21	7440-42-8	CAS-RN	DTXSID3023922	Boron	-	-	-	Υ
22	7440-29-1	CAS-RN	DTXSID6049800	Thorium	-	-	-	Υ

List of Chemicals





EPA activities around PFAS chemicals

https://www.epa.gov/pfas



Per- and Polyfluoroalkyl Substances (PFAS)



CONTACT US

SHARE









"The National Leadership
Summit on PFAS provided
an unprecedented
opportunity for stakeholders
to share vital information
and best practices regarding
PFAS." -

Former Administrator Pruitt

- Community Events
- Infographic

Basic Information

- What are PFAS?
- Why are PFAS important?
- How people are exposed?

EPA Actions to Address PFAS

- EPA actions
- National leadership summit and engagement

Tools and Resources

- EPA data and tools
- State information
- Site-specific resources

The OECD List of PFAS

http://www.oecd.org/chemicalsafety/portal-perfluorinated-chemicals/







HOME



The OEGD releases a new list of PFASs

The OECD releases a new list of Per- and Polyfluoroalkyl Substances (PFASs) based on a comprehensive analysis of information available in the public domain. In total, 4730 PFAS-related CAS numbers have been identified and categorised in this study, including several new groups of PFASs that fulfil the common definition of PFASs (i.e. they contain at least one perfluoroalkyl moiety) but have not yet been commonly regarded as PFASs.

This work has been conducted under the OECD/UN Environment Global PFC Group in support of the Strategic Approach to International Chemicals Management (SAICM) and shifting to safer alternatives for PFASs.

The New Comprehensive Global Database of Per- and Polyfluoroalkyl Substances (PFASs) comes with a methodology report also detailing the major findings with respect to the total numbers and types of PFASs identified, the limitations, gaps and challenges identified in the development of the new list, and opportunities for improving the future understanding of PFASs production, use on the global market, and presence in the environment, biota, and other matrices.





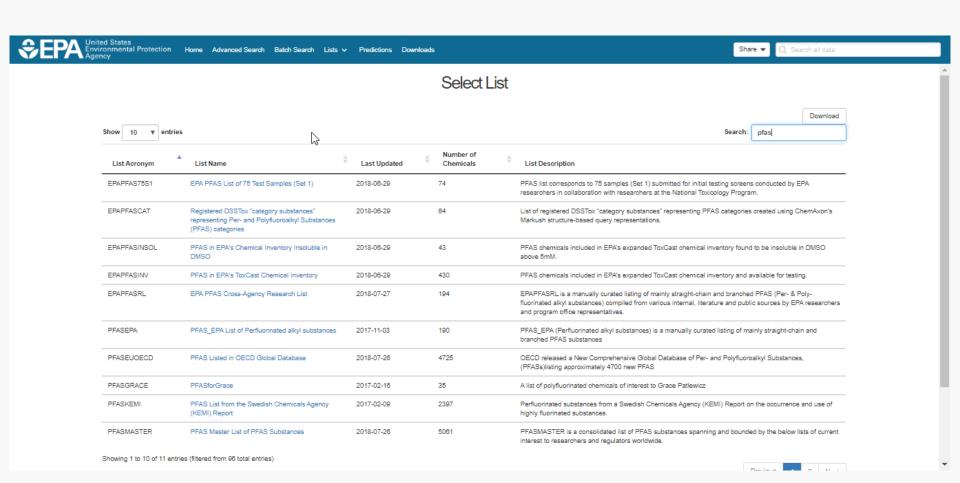




11 PFAS Lists

http://comptox-prod.epa.gov/dashboard/chemical_lists





PFAS Categories in Development



Registered DSSTox "category substances" representing Per- and Polyfluoroalkyl Substances (PFAS) categories

Search EPAPFASCAT Chemicals	Q
Substring search	

List Details

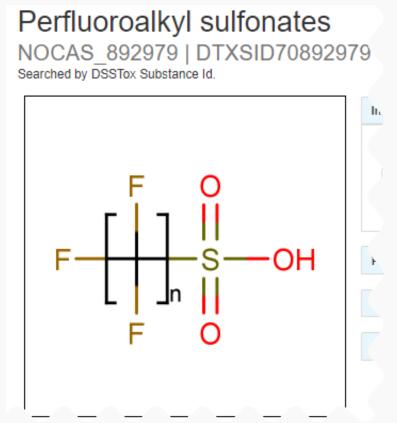
Description: List of registered DSSTox "category substances" representing Per- and Polyfluoroalkyl Substances (PFAS) categories created using ChemAxon's Markush structure-based query representations. Markush categories can be broad and inclusive of more specific categories, or can represent a unique category not overlapping with other registered categories. Each PFAS category registered with a unique DTXSID is considered a generalized substance or "parent ID" that can be associated with one or many "child IDs" (i.e. many parent-child mappings) within the full DSSTox database. These category DTXSIDs can be used to search and retrieve all currently registered DSSTox substances within the category group, and offer an objective, transparent and reproducible structure-based means of defining a category of chemicals. This list and the corresponding category mappings is undergoing continuous curation and expansion.

Number of Chemicals: 64

Markush Chemicals



PFOS is a member of linear perfluoroalkyl sulfonates



Collaborative Data Curation



 Mapping between our data (and websites) has resulted in collaborative data curation

 Collaboration with Emma Schymanski re. the NORMAN Suspects Exchange https://www.norman-network.com/?q=node/236

Multiple NORMAN lists now mapped

NORMAN Suspect Exchange





Network of reference laboratories, research centres and related organisations for monitoring of emerging environmental substances

NORMAN Network Working Groups NORMAN Bulletin Success Stories **Publications** NORMAN GA meetings Membership Job opportunities Members' Area

Home Emerging Substances

DATABASES

Menu

- Topics and Activities
- Workshops and Events
- QA/QC Issues
- Glossary
- Useful links
- Members' Area



http://www.norman-network.com/?q=node/236

NORMAN Suspect List Exchange

In September 2014, NORMAN members expressed the need to exchange various lists of substances to improve their suspect screening efforts. This website was established aspart of the 2015 Joint Programme of Activities as a central access point for NORMAN members (and others) to find suspect lists relevant for their environmental monitoring question. All suspect lists currently available are compiled in the table below and on the US EPA CompTox Chemistry Dashboard (website, downloads, chemical lists).

The "Link to full list" column below contains an excel or comma-separated file (csv) with all available information, e.g. as provided as supporting information for the publication, while the third column provides a list of the structures as InChlKeys only, which allows suspect searching using MetFrag or other workflows. The fourth column contains references for the data: please cite these references if you use the respective datasets.

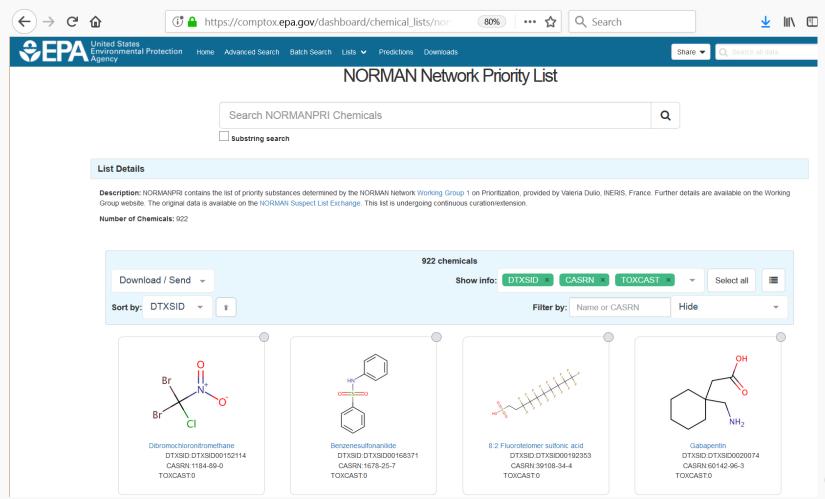
Recent Suspect Exchange and Dashboard presentations/publications include: ICCE Oslo 2017: NORMAN Suspects meet the Dashboard and NORMAN MassBank and Suspect Exchange; SETAC Mixtures Denver: Identifying Complex Mixtures with Cheminformatics and HR-MS; ACS Fall 2017: Markush Enumeration for UVCBs and a viewpoint article.

No.	Abbreviation	Description	Link to full list	Link to InChlKey list	References
	SUSDAT	Merged NORMAN Suspect List: SusDat	Interactive Data table (updating)	MS-ready InChlKeys (1/03/2018)	A merged list of >40,000structures from suspect lists. See interactive version. Compiled by Reza Aalizadeh, University of Athens, including RTI and toxicity values, support by Nikiforos Alygizakis, El. Work in progress please report any issues!
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 (11/04/2017)	www.massbank.eu Stravs et al. 2013. DOI: 10.1002/jms.3131
S2	STOFFIDENT	HSWT/LfU STOFF- IDENT Database of Water-Relevant Substances	STOFF-IDENT Contents (6/09/2017) CompTox STOFF-IDENT List Further curation in progress	STOFF-IDENT InChlKeys (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).
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 (31/10	Schymanski <i>et al.</i> 2015. DOI: 10.1007/s00216-015-8681-7

Example: NORMAN Priority List

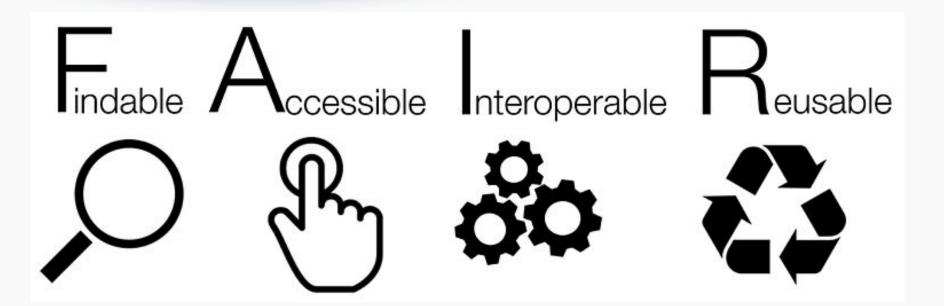


				1	
S15	NORMANPRI	NORMAN Priority List	NORMAN Priority CSV (13/7/2017)	NORMAN Priority	Priority substances from NORMAN WG-1 (Prioritisation),
			CompTox NORMAN Priority List	InChlKeys (16/05/2017)	provided by Valeria Dulio.
			Further curation in progress		



Our support for FAIR Data





 We're definitely not there yet...but are making progress at speed...

How we are serving FAIR



- Our data are licensed as public domain data
 - available from downloads page
 - registered to Figshare
 - SQL data dumps
- Collection of web services for old dashboards are available – API is being fully revamped

API in development Prototype services available



https://comptox.epa.gov/dashboard/web-test/WS?smiles=CCO&method=hc

```
JSON
         Raw Data
                    Headers
Save Copy
 uuid:
                           "55547f4f-f966-48e8-b831-a0d217998064"
 predictionTime:
                           1520539090089
 software:
                           "T.E.S.T (Toxicity Estimation Software Tool)"
 softwareVersion:
                           "5.01"
                           "25°C"
 condition:
                           "Water solubility at 25°C"
 endpoint:
 method:
                           "Hierarchical clustering"
▼ predictions:
  ₹0:
                           "C 1520539090089"
       id:
       smiles:
                           "0CC"
       expValMolarLog:
                           "-1.337"
                           "1001180.703"
       expValMass:
                           "-1.338"
       predValMolarLog:
       predValMass:
                           "1002625.241"
                           "-Log10(mo1/L)"
       molarLogUnits:
                           "mg/L"
       massUnits:
```

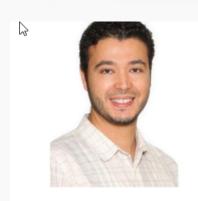
How we are serving FAIR



- Our data are licensed as public domain data
 - available from downloads page
 - registered to Figshare
 - SQL data dumps
- Collection of web services for old dashboards are available – API is being fully revamped
- Models are published to Github

OPERA Models on Github https://github.com/kmansouri





Kamel Mansouri

kmansouri

- Computational Chemistry/Toxicology -Cheminformatics and data mining -QSAR/QSPR and ADME-Tox properties modeling [orcid:0000-0002-6426-8036]

Block or report user

RTP, NC, USA

യ https://www.linkedin.com/in/ka...

Overview Repositories 3 Stars 3

Followers 5 Following 4

Popular repositories

OPERA

Command line application providing QSAR models predictions as well as applicability domain and accuracy assessment for physicochemical properties and environmental fate endpoints.

★3 ¥1

QSAR-ready

Standardization workflow for QSAR-ready chemical structures pretreatment.

MS-ready

Standardization workflow for MS-ready chemical structures pretreatment.

12 contributions in the last year



How we are serving FAIR



- Our data are licensed as public domain data
 - available from downloads page
 - registered to Figshare
 - SQL data dumps
- Collection of web services for old dashboards are available – API is being fully revamped
- Models are published to Github
- DTXSIDs accepted on WikiData

DTXSIDs on WikiData





Main page Community portal Project chat Create a new item

Recent changes Random item

Query Service Nearby

Help Donate

Print/export

Daumland on DDF

Property Discussion

文 English L Not

Read View history

Search Wi

DSSTOX substance identifier (P3117)

DSSTox substance identifier used in the Environmental Protection Agency CompTox Dashboard DTXSID

▼ In more languages Configure

Language	Label	Description	Also known as
English	DSSTOX substance identifier	DSSTox substance identifier used in the Environmental Protection Agency CompTox Dashboard	DTXSID
German	DSSTOX-Identifikator	No description defined	DTXSID
French	identifiant DSSTOX	identifiant DSSTox d'une substance utilisé par l'agence de protection de l'environnement américaine	DTXSID

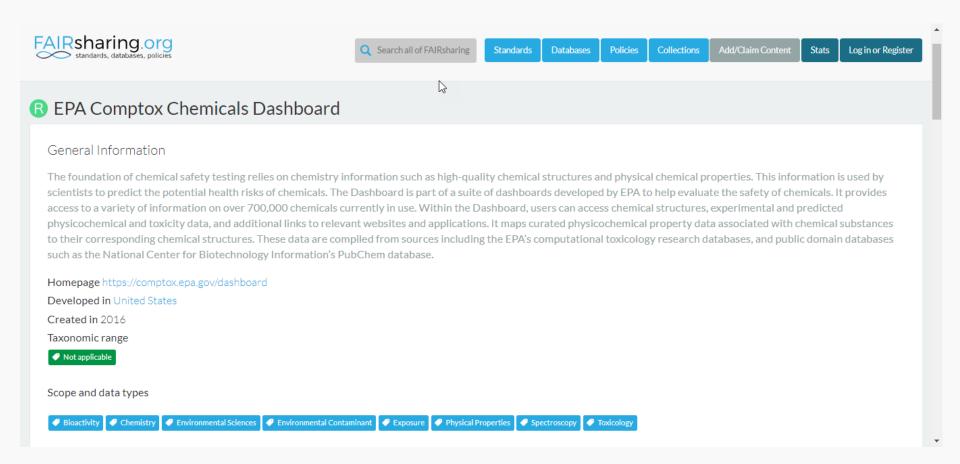
How we are serving FAIR



- Our data are licensed as public domain data
 - available from downloads page
 - registered to Figshare
 - SQL data dumps
- Collection of web services for old dashboards are available – API is being fully revamped
- Models are published to Github
- DTXSIDs accepted on Wikidata
- In discussions now re. generation of RDF
- FAIR page will be updated with progress

FAIRsharing.org page https://fairsharing.org/FAIRsharing.tfj7gt





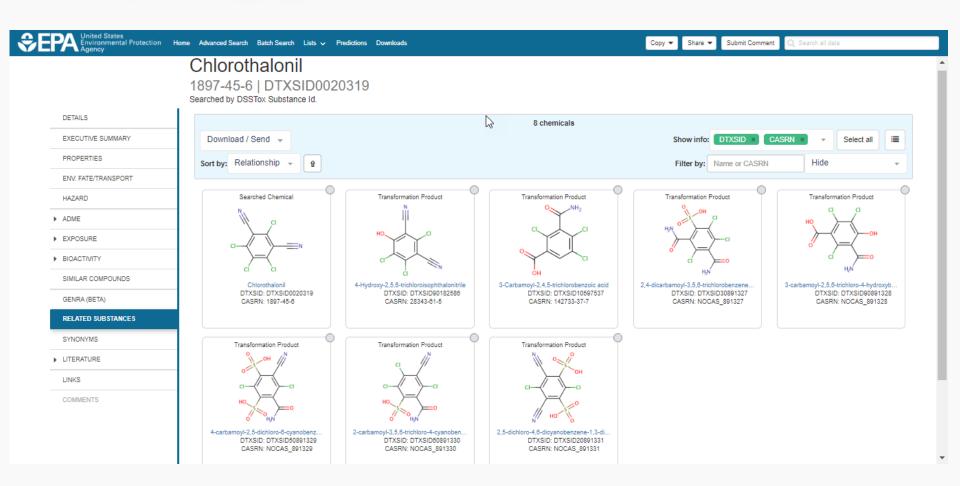
In Progress: "InvitroDB_v3"



- The last public release of ToxCast data (invitroDB_v2) was in 3rd Quarter of 2015
- Next release of invitroDB_v3 is Fall 2018
- Data includes new assays, new chemicals, new pipelining, results of data curation
- Data will also release via the Dashboard
- Data available at https://www.epa.gov/chemical-research/exploring-toxcast-data-downloadable-data

In Progress: Metabolites extracted from literature and databases

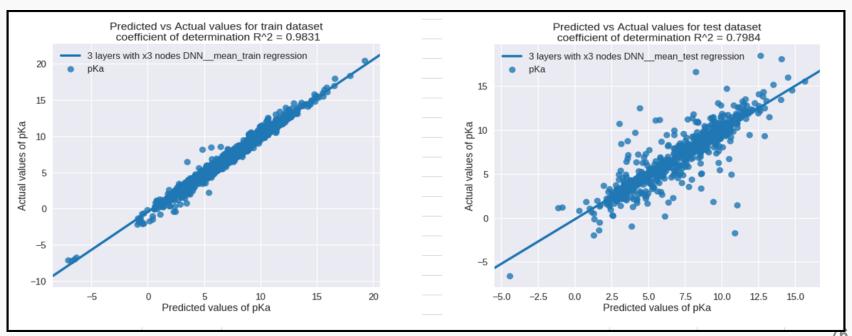




In Progress : pKa Prediction Model



pKa prediction models based on Open
 Data Set of 8000 chemicals – acidic, basic
 and amphoteric chemicals



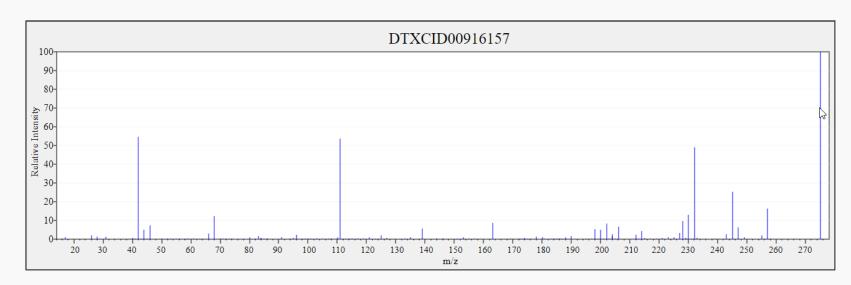
In Progress: 700k Predicted MS

http://cfmid.wishartlab.com/





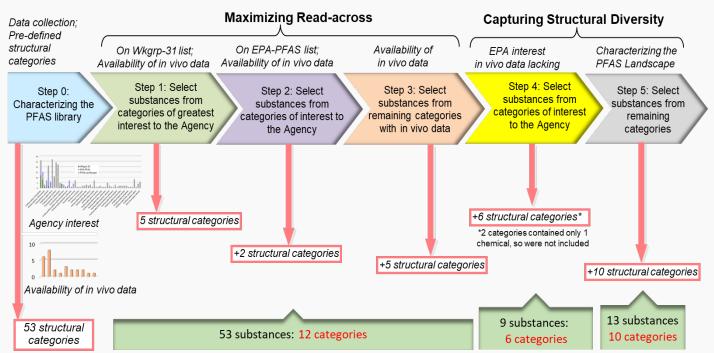
- MS/MS spectra prediction for ESI+, ESI-, and EI
- Predictions generated and stored for >700,000 structures, to be accessible via Dashboard



In progress: PFAS Chemical Library



- Development of a high-throughput screening library and collection of physical samples (~400)
- 75 PFAS chemicals for screening based on categories, diversity, exposure considerations, procurability and testability, availability of existing toxicity data



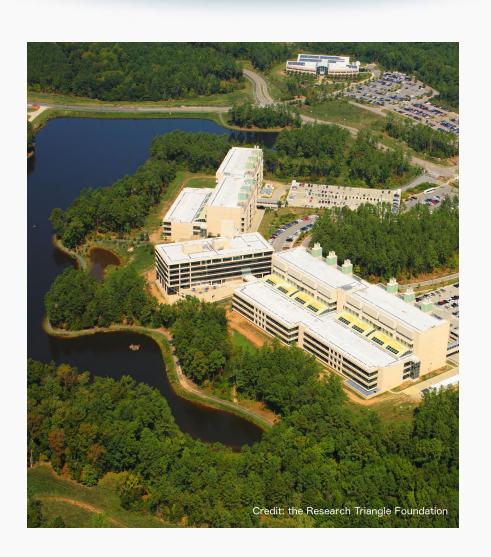
Conclusion



- Transparent access to data supporting computational toxicology – file downloads, SQL data dumps and web services
- CompTox Chemicals Dashboard provides access to data for ~765,000 chemicals
- Our publications provide chemistry data in usable formats – Excel, SDF (V2000/V3000)
- Web Services developing to serve API development
- Next release: Fall 2018 with InvitroDB v3 data

Acknowledgements





EPA-RTP

- An enormous team of contributors from NCCT
- Collaborators from multiple EPA labs and centers
 - Kamel Mansouri (OPERA)
 - Todd Martin (TEST)
 - Valery Tkachenko (TEST)
- Emma Schymanski and the NORMAN Network (LCSB)

Contact

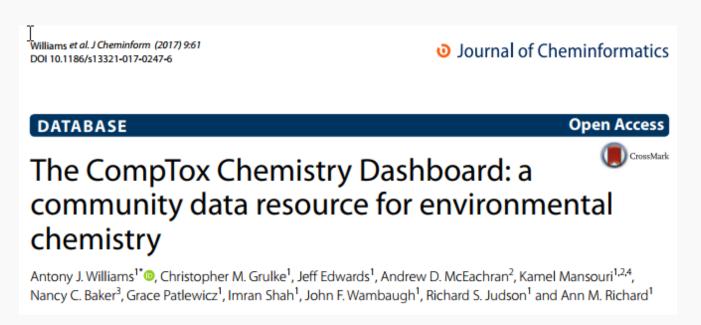


Antony Williams

NCCT, US EPA Office of Research and Development,

Williams.Antony@epa.gov

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



https://doi.org/10.1186/s13321-017-0247-6