

# The advantages, and cautions, of applying AI to the development of prediction models

Antony Williams

AAAS Conference, Denver CO. February 2024

The views expressed in this presentation are those of the author and do not necessarily reflect the views or policies of the U.S. EPA

### US-EPA Office of Research and Development

- Center for Computational Toxicology and Exposure (CCTE)
- Some of our primary efforts in AI are in modeling (QSAR/QSPR QSUR) based on measurement, harvesting, assembly and curation of streams of data
- We develop and deliver to the community
  - curated chemistry data streams to support our applications and models
  - prediction models, web-based applications and data streams to support others



• We have benefited from AI modeling approaches for decades

Approaches continue to advance with the promise of improved performance

 Many challenges remain the same – quality of data, data overfitting, and interpretability of the models

• The latest form of modernization is "Large Language Models"



 There are many tools developed by our cheminformatics team and across other centers in EPA. I will represent ours only...

• We have production level public-facing tools, proof-of-concept public-facing tools, and many tools in development...

 We focus on FAIR data releasing it to the community and making it available on Public APIs

### Example Problems We Apply Models to



- Property prediction e.g., water solubility, vapor pressure
- Fate and Transport e.g., bioaccumulation, bioconcentration
- Bioactivity e.g., endocrine disruption, steroidogenesis
- Analytical Chemistry e.g., what techniques are chemicals amenable to - GC/MS? LC/MS? +ve/-ve ion?

### **OECD** Principles for Modeling

https://www.oecd.org/chemicalsafety/risk-assessment/37849783.pdf



- To facilitate the consideration of a (Q)SAR model for regulatory purposes, it should be associated with the following information:
  - 1) a defined endpoint
  - 2) an unambiguous algorithm
  - 3) a defined domain of applicability
  - 4) appropriate measures of goodness-of-fit, robustness and predictivity
  - 5) a mechanistic interpretation, if possible
- These principles have been around a long time...

These principles were agreed by OECD member countries at the 37<sup>th</sup> Joint Meeting of the Chemicals Committee and Working Party on Chemicals, Pesticides and Biotechnology in November 2004.



### Data assembly and curation



- ~10,200 unique chemical structures with experimental data
- Great care is needed in data assembly for modeling
- For water solubility large collections based on standard methods provide threshold measurements >X, <Y</li>
- The same data is shared between public platforms





- Many Descriptors to choose: commercial and open source
- We use Padel, Mordred and TEST descriptors (open)
- Example: <a href="http://www.yapcwsoft.com/dd/padeldescriptor/">http://www.yapcwsoft.com/dd/padeldescriptor/</a>

## PaDEL-Descriptor

### Description

A software to calculate molecular descriptors and fingerprints. The software currently calculates 1875 descriptors (1444 1D, 2D descriptors and 431 3D descriptors) and 12 types of fingerprints (total 16092 bits). The descriptors and fingerprints are calculated using The Chemistry Development Kit with additional descriptors and fingerprints such as atom type electrotopological state descriptors, Crippen's logP and MR, extended topochemical atom (ETA) descriptors, McGowan volume, molecular linear free energy relation descriptors, ring counts, count of chemical substructures identified by Laggner, and binary fingerprints and count of chemical substructures identified by Laggner.

# Feature Selection and Variables can help mechanistic understanding





Without Feature Selection – 427 variables

With Feature Selection – 19 variables

### Descriptors



Descriptor	Definition
XLOGP	Wang octanol water partition coefficient
XLOGP2	Wang octanol water partition coefficient squared
BEHm3	Highest eigenvalue n. 3 of Burden matrix / weighted by atomic masses
eim	Mean information content on the edge magnitude
Qs	Molecular and group <b>polarity</b> index
BEHm4	Highest eigenvalue n. 4 of Burden matrix / weighted by atomic masses
BEHm2	Highest eigenvalue n. 2 of Burden matrix / weighted by atomic masses
BEHm7	Highest eigenvalue n. 7 of Burden matrix / weighted by atomic masses
BEHm6	Highest eigenvalue n. 6 of Burden matrix / weighted by atomic masses
MATS1v	Moran autocorrelation - lag 1 / weighted by atomic van der Waals volumes
MATS1p	Moran autocorrelation - lag 1 / weighted by atomic polarizabilities
ATS4m	Broto-Moreau autocorrelation of a topological structure - lag 4 / weighted by atomic masses
BEHm5	Highest eigenvalue n. 5 of Burden matrix / weighted by atomic masses
Hmax	Maximum hydrogen E-State value in molecule.
pilD	Conventional bond order id number
ieadjem	Mean information content on the edge adjacency equality
Mv	Mean atomic van der Waals volume (scaled on Carbon atom)
SHHBd	Sum of E-State indices for hydrogen bond donors
ide	Total information content on the distance equality



- Many of our tools are publicly available
  - CompTox Chemicals Dashboard (https://comptox.epa.gov/dashboard)
    - Predicted properties (ACD/Labs, OPERA, TEST)
    - Bioactivity models: Estrogen receptor, Androgen receptor
  - Proof-of-Concept cheminformatics modules (<u>https://www.epa.gov/chemical-research/cheminformatics</u>)
    - Hazard Profiling of chemicals
    - Single and batch prediction of physprop and toxicity endpoints
- Chemical curation is into the DSSTox database

### Assembling data is easy. Curation is hard

https://pubs.acs.org/doi/10.1021/acs.jcim.2c00268



 It is very easy to harvest and download massive amounts of data. FAIRness has expanded access...

#### RETURN TO ISSUE < PREV CHEMICAL INFORMATION NEXT >

### CAS Common Chemistry in 2021: Expanding Access to Trusted Chemical Information for the Scientific Community

Andrea Jacobs\*, Dustin Williams, Katherine Hickey, Nathan Patrick, Antony J. Williams, Stuart Chalk, Leah McEwen, Egon Willighagen, Martin Walker, Evan Bolton, Gabriel Sinclair, and Adam Sanford

 Open API and downloadable dataset – contributing CASRNs, Names and Structures to Open Chemistry

### Stoichiometry is important

CAS



### • SIMPLE example...1 to 3 stoichiometry

#### Benzoic acid, praseodymium(3+) salt (3:1)

**Common Chemistry** 



### **Other Names and Identifiers**

InChI=1S/C7H6O2.Pr/c8-7(9)6-4-2-1-3-5-6;/h1-5H,(H,8,9);

InChIKey=ZUSLIYUUHQQOHU-UHFFFAOYSA-N

SMILES

C(O)(=O)C1=CC=CC=C1.[Pr]

Canonical SMILES [Pr].O=C(O)C=1C=CC=CC1

1000s of structures with bad stoichiometry into the wild





- A publicly accessible website delivering:
  - 1.2M chemicals with related property data
  - Related substances: transformation products, mono/polymer
  - Experimental/predicted physicochemical property data
  - Experimental Human and Ecological hazard data
  - Integration to "biological assay data" (ToxCast/Tox21)
  - Information regarding chemicals in consumer products
  - Links to other agency websites and public data resources
  - "Batch searching" for tens to thousands of chemicals

### Curating Chemistry into the DSSTox Database



Computational Toxicology Volume 12, November 2019, 100096



# EPA's DSSTox database: History of development of a curated chemistry resource supporting computational toxicology research

<u>Christopher M. Grulke</u><sup>a</sup>, <u>Antony J. Williams</u><sup>a</sup>, <u>Inthirany Thillanadarajah</u><sup>b</sup>, <u>Ann M. Richard</u><sup>a</sup> *A* ⊠



mental Protection

### CompTox Chemicals Dashboard

### https://comptox.epa.gov/dashboard



CompTox Chemicals Dashb	oard Home Search ▼ Lists ▼ About ▼ Tools ▼			Sul	bmit (	Comm	nents
Mac	CompTox Chemicals Dashboard Search 1,200,059 Chemicals						
Same.	Chemicals Products/Use Categories Assay/Gene						
	Perfluoro       Perfluoro-10-dodecene-1-sulfonic acid       Putrolanagenega	^	٩				L
DARK	perfluoro-10-dodecene-1-sulfonic acid         DTXSIDB01032598         Perfluoro-10-oxoundecane-1-sulfonic acid         DTXSIDB01032379         Perfluoro-10-ixoundecane-1-sulfonic acid         DTXSIDB01032379         Perfluoro-10-(pentafluoro-lambda6-sulfanyl)decanoic acid         DTXSIDB01035809         Perfluoro-1-(1308)octanesulfonamide         Perfluoro-1-(1308)octanesulfonamide		٩	ALL ALL ALL			

### 1 of ~1.2M Chemical Pages Experimental/Predicted Properties



CompTox Chemicals Dashboard Home Search -	Lists - About - Tools -	Submit Comments Search all data	<b>∨</b> Q
Perfluc       +++++++++       1763-23       ©	Prooctanesulfonic acid B-1   DTXSID3031864 DTXSID		
Details     Chemical Details       Executive Summary     Properties       Env. Fate/Transport     F       Hazard     F       Safety > GHS Data     F	Wikipedia           Perfluorooctanesulfonic acid (PFOS) (conjugate base perfluorooctanesulfonate functional group and thus a perfluorosulfonic acid. It is an anthropogenic (man-m Scotchgard, a fabric protector made by 3M, and related stain repellents. The acro Read more	) is a chemical compound having an eight-carbon fluorocarbon chain and a sulfonic acid nade) fluorosurfactant, now regarded as a global pollutant. PFOS was the key ingredient in onym "PFOS" refers to the parent sulfonic acid and to various salts	•
ADME > IVIVE Exposure Bioactivity Similar Compounds GenRA Related Substances	Quality Control Notes         Intrinsic Properties         Image: Molecular Formula: CgHF1703S       Mol. FILE       Q FIND ALL CF         Average Mass: 500.13 g/mol       Image: Source File       Monoisotopic Mass: 499.937494 g/mol	* HEMICALS	
Synonyms Literature Links Comments	Structural Identifiers Linked Substances Presence in Lists	* * *	

•

### QSAR Modeled Data are available



- We build models then apply then to our curated datasets for release, PLUS deliver the models for realtime use
- Data can be harvested one chemical at a time or 1000s using the batch search

(https://pubs.acs.org/doi/10.1021/acs.jcim.0c01273)



**Enabling High-Throughput Searches for Multiple Chemical** Data Using the U.S.-EPA CompTox Chemicals Dashboard

Charles N. Lowe\* and Antony J. Williams\*

✓ Cite this: J. Chem. Inf. Model. 2021, 61, 2, 565–570 Publication Date: January 22, 2021 ~ 1332 https://doi.org/10.1021/acs.jcim.0c01273

Article Views

Altmetric

2

Citations

32

### Modeled values are not enough!





#### **Bisphenol A** 80-05-7 | DTXSID7020182

Searched by Expert Validated Synonym.

#### **Properties: Summary**

Chemical Details 👗	Properties: Summ	ary					
Executive Summary	Summary	× ۵	Search Chemical Properties				
Physchem Prop.				•			
Env. Fate/Transport	LEXPORT -			Summary			
Hazard Data	Property ↓↑ Ξ	$\equiv$ Experimental average $\downarrow$	$\equiv ig $ Predicted average $\downarrow\uparrow$	$\equiv ig $ Experimental median $\downarrow\uparrow$	$\equiv$   Predicted median $\downarrow\uparrow$	$\equiv$ Experimental range $\downarrow\uparrow$	=
Safety > GHS Data	Water Solubility	8.55e-4 (3)	7.81e-4 (4)	5.26e-4	7.49e-4	5.25e-4 to 1.51e-3	
ADME > IVIVE	LogKow: Octanol-Water	3.32 (1)	3.50 (4)	3.32	3.53	3.32	
	Boiling Point	200 (1)	367 (4)	200	362	200	
Exposure	Melting Point	155 (7)	137 (3)	156	132	153 to 156	
Bioactivity	Polarizability	-	27.0 (1)	-	27.0	-	
	Henry's Law	-	1.26e-7 (1)	-	1.26e-7	-	
GenRA	ReadyBiodeg	-	0.00 (1)	-	0.00	-	
ACToR	Flash Point	-	190 (2)	-	190	-	
<del>.</del>	Molar Refractivity		68.2 (1)		68.2		_

Where is all the calculation detail? Are predictions in applicability domain etc?



- For OPERA and TEST models we have all the details
  - OPERA https://jcheminf.biomedcentral.com/articles/10.1186/s13321-018-0263-1

RESEARCH ARTICLE	Open Access
OPERA models for predicting physicochemical properties and environr fate endpoints	orossMark) nental
Kamel Mansouri <sup>1,2,3*</sup> Chris M. Grulke <sup>1</sup> . Richard S. Judson <sup>1</sup> and Antony J. Williams <sup>1</sup>	

- TEST <u>https://www.epa.gov/comptox-tools/toxicity-estimation-software-tool-test</u>



### **OPERA Model Details**



SEPA

United States Environmental Protection Agency

### **OPERA Model Details**





### **OPERA Model Details**



#### Nearest Neighbors from the Training Set



### Why this detail is required



### • Predicted Fish Biotrans. Half-Life (Km) of PFOS is 2.7 days

#### OPERA Model Calculation Details: Fish Biotrans. Half-Life (Km)

#### 🕒 PRINT PDF

Chemical Identifiers	Model Results
Preferred Name: Perfluorooctanesulfonic acid DTXSID: DTXSID3031864 DTXCID: DTXCID1011864 CASRN: 1763-23-1	Predicted value: 2.69 days Global applicability domain: Outside Local applicability domain index: 0.174 Confidence level: 0.145 Opera Version: OPERA 2.6

### Why this detail is required



#### Model Performance KM data KM model 70 Training set Test set 60 50 40 30 20 Training set 10 Test set -1.5 h -2 -1.5 -0.5 0.5 1.5 2 2.5 -1 0 0 LogKM observed -1.5 -0.5 0.5 2 2.5 -2 -1 0 1 1.5 3 📥 QMRF

#### Weighted KNN model

5-fold CV (75%)		Training	(75%)	Test	t (25%)
Q2	RMSE	RMSE	R2	R2	RMSE
0.830	0.490	0.500	0.820	0.730	0.620

### QMRF details

#### https://comptox.epa.gov/dashboard-api/ccdapp1/qmrfdata/file/by-modelid/28



# OMRF

#### QMRF identifier (JRC Inventory):Q17-66-0019 **OMRF Title:OPERA-model** for biotransformation rate constant Printing Date:Oct 17, 2017

#### 1.OSAR identifier

#### **1.1.QSAR identifier (title): OPERA-model** for biotransformation

#### **1.2.Other related models:** No related models

#### **1.3.Software coding the model: OPERA V1.5**

OPERA (OPEn (quantitative) str open source command line appl physicochemical properties and descriptors. It is available for do license.

Kamel Mansouri (mansourikame https://github.com/kmansouri/OI

#### **3.7.Endpoint data quality and variability:**

The original data collected from the PHYSPROP database (631 chemicals) have undergone a se 4.4.Descriptor selection: the chemical structures and remove duplicates, erroneous entries. This procedure also included ensure only good quality data is used for the de model (548 chemicals).

- Then, QSAR-ready structures were generated by standardizing all chemical structures and remov and metallo-organic chemicals (541 chemicals). workflows that were developed for the purpose standardization of the data are available in the r Section 2.7].
- The curated outlier-free experimental data (541 chemicals) was divided into training and validati machine learning and modeling steps.

PaDEL software was used to calculate 1440

molecular descriptors. A first filter was applied in order to remove descriptors with missing values, constant and near constant (standard deviation of 0.25 as a threshold) and highly correlated descriptors (96% as a threshold). The remaining 837 descriptors were used in a feature selection procedure to select a minimum number of variables encoding the most relevant structural information to the modeled endpoint. This step consisted of coupling Genetic Algorithms (GA) with the weighted kNN algorithm and was applied in 5 fold cross validation on the training set (405 chemicals). This procedure was run for 200 consecutive independent runs maximizing  $Q^2$  in cross-validation and minimizing the number of descriptors. The number of k neighbors is optimized within the range of 3 to 7. The descriptors were then ranked based on their frequency of selection during the GA runs. The best model showed an optimal compromise between the simplicity (minimum number of descriptors) and performance (Q<sup>2</sup> in cross-validation) to ensure transparency and facilitate the mechanistic interpretation as required by OECD principles 2 and 5. More details in paper [ref2 Section 2.7].

### Why this detail is required



#### Nearest Neighbors from the Training Set



### **TEST Prediction Reports**



Predicted Water solubility at 25 ŰC for 80-05-7 from Consensus method

Prediction results				
Endpoint	Experimental value (CAS= 80-05-7) Source: <u>EPI Suite v 4.00</u>	Predicted value <sup>a</sup>		
Water solubility at 25 ŰC -Log10(mol/L)	3.28	2.90		
Water solubility at 25 ŰC mg/L	120.10	284.38		

<sup>a</sup>Note: the test chemical was present in the training set. The prediction *does not* represent an external prediction.



Predictions for the test chemical and for the most similar chemicals in the e

If the predicted value matches the experimental values for similar chemicals in the test set (and t



Chemicals	MAE*		
Entire set	0.58		
Similarity coefficient $\ge 0.5$	0.42		
Mean absolute error in -Log10(mol/L)			

CAS	Structure	Similarity Coefficient	Experimental value -Log10(mol/L)	Predicted value -Log10(mol/L)
80-05-7 (test chemical)	s d d		3.28	2.90
<u>77-21-4</u>	, (fO	0.72	2.34	2.60
<u>103387-44-6</u>	≺.Ĵ.			3.09
94-77-9	Û.,	0.69	2.72	3.19
<u>491-59-8</u>		0.61	3.06	4.03
<u>19578-81-5</u>		0.60	4.64	4.23
110048-81-2	f.	0.60	3.43	3.40
<u>15307-86-5</u>		0.59	5.10	4.47
<u>101-77-9</u>	.m.	0.58	2.30	2.97

### Access to Real Time Predictions





Toxicological properties
 96 hour fathead minnow LC50
 48 hour D. magna LC50
 48 hour T. pyriformis IGC50
 Oral rat LD50
 Bioconcentration factor
 Developmental toxicity
 Ames mutagenicity
 Estrogen Receptor RBA
 Estrogen Receptor Binding

Physical properties
 Normal boiling point
 Melting point
 Flash point
 Vapor pressure
 Density
 Surface tension
 Thermal conductivity
 Viscosity
 Water solubility

### Multiple modeling approaches plus Consensus



Property =	Experimental Value		$\equiv$ Hierarchical clustering	$\equiv$ Single model	$\equiv$ Group contribution	
96 hour fathead minnow LC50	4.158 -Log10(mol/L) 14.993 mg/L	4.001 -Log10(mol/L) 21.524 mg/L	4.067 -Log10(mol/L) 18.489 mg/L	3.997 -Log10(mol/L) 21.697 mg/L	3.876 -Log10(mol/L) 28.721 mg/L	4.064 -Log10(mol/L)
48 hour D. magna LC50	3.601 -Log10(mol/L) 54.062 mg/L	3.854 -Log10(mol/L) 30.189 mg/L	3.725 -Log10(mol/L) 40.679 mg/L	4.354 -Log10(mol/L) 9.555 mg/L	3.817 -Log10(mol/L) 32.852 mg/L	3.521 -Log10(mol/L) 65.046 mg/L
48 hour T. pyriformis IGC50			3.525 -Log10(mol/L) 64.465 mg/L			
Oral rat LD50	2.506 -Log10(mol/kg) 672.807 mg/kg	2.197 -Log10(mol/kg) 1370.798 mg/kg	2.295 -Log10(mol/kg) 1093.260 mg/kg			2.099 -Log10(mol/kg) 1718.794 mg/kg
Bioconcentration factor	0.557 Log10 3.606	0.950 Log10 8.910	0.837 Log10 6.873	0.985 Log10 9.650	1.004 Log10 10.102	0.974 Log10 9.409
Developmental toxicity		true	true	true		false
Ames mutagenicity	false	false	false			false
Estrogen Receptor RBA						
Estrogen Receptor Binding	false	false	false	false	false	
Normal boiling point		312.2 °C	310.7 °C		341.5 °C	284.3 °⊂
Melting point	173.0 °C	163.6 °C	180.6 °C		106.2 °C	204.0 °⊂
•						• •
Rows: 18		1	otal Rows: 18			

### Coming Soon: Excel report for models for each data set



- Cover sheet with model metadata
- Training and test set statistics

- Training and test set statistics
- Prediction results for each method



### Where do we use predictions like this?



- Models are used in many places in our computational toxicology research
- They are used in the analytical labs to help guide nontargeted analysis



### Where do we use predictions like this?

- Models are used in many places in our computational toxicology research
- They are used in the analytical labs to help guide nontargeted analysis
- By stakeholders for Hazard profiling of chemicals



vironmental Protection

### Where do we use predictions like this?

- Models are used in many places in our computational toxicology research
- They are used in the analytical labs to help guide nontargeted analysis
- By stakeholders for Hazard profiling of chemicals
- Predictions for breakdown products in the environment



nmental Protection

### Where can our tools be applied



- Emergency Response utility is obvious...
- Consider East Palestine



https://www.cleveland19.com/2023/ 02/14/ntsb-announces-preliminarymalfunction-that-caused-eastpalestine-train-derailment/

POLYPROPYLENE POLYETHYLENE Residue lube oil VINYL CHLORIDE DIPROPYLENE GLYCOL **PROPYLENE GLYCOL** DIETHYLENE GLYCOL COMBUSTIBLE LIQ., NOS (ETHYLENE GLYCOL MONOBUTYL ETHER) SEMOLINA COMBUSTIBLE LIQ., NOS (ETHYLHEXYL ACRYLATE) POLYVINYL PETROLEUM LUBEOIL POLYPROPYL GLYCOL **ISOBUTYLENE** BUTYL ACRYLATES, STABILIZED PETRO OIL, NEC ADDITIVES, FUEL BALLS, CTN, M EDCL SHEET STEEL VEGTABLE, FROZEN BENZENE PARAFFIN WAX FLAKES, POWDER HYDRAULIC CEMENT AUTOS PASSENGER MALT LIQUORS

### Hazard Comparison Profiling



	result	
BENZENE VINYL CHLORIDE DIPROPYLENE GLYCOL PROPYLENE GLYCOL DIETHYLENE GLYCOL ETHYLENE GLYCOL MONOBUTYL ETHER POLYVINYL		
POLYPROPYL GLYCOL ISOBUTYLENE BUTYL ACRYLATE	Search by identifiers Retrieve Comptox List Sea	arch by structure Search result
	Chemical(s) found by 14 identifier(s)	12 / 15 🕑 🔲 🖽 💥
	Benzene         Vinyl chloride         Dipropylene           71-43-2         75-01-4         25265-7           BENZENE         VINYL CHLORIDE         DIPROPYLENE	e glycol 1.2-Propylene glycol 1-8 57-55-6 E GLYCOL PROPYLENE GLYCOL Diethylene glycol 111-46-6 DIETHYLENE GLYCOL DIETHYLENE GLYCOL
	Isobutene Butyl acrylate Ethylhexyl a 115-11-7 141-32-2 1322-13	Acrylate Hydrocarbon waxes Polypropylene Polyethylene AS low 3-0 8002-74-2 9003-07-0 9002-88-4

### Hazard Comparison Profiling



Chemicals: 12 Toxicity: VH - Very High H - High M - Medium L - Low I - Inconclusive N/A - Not Applicable Authority: Authoritative O Screening QSAR Model O														Nodel 🛈						
Chomodalo. 12	Human Health Effects													Ecoto	oxicity		Fate			
<ul> <li>Skipped (0)</li> <li>Unlikely (0)</li> <li>Filters (0)</li> <li>Sorting (0)</li> <li>Structure CAS Name</li> </ul>	Acute M	lammaliar uoite	n Toxicity Derman	Carcinogenicity	Genotoxicity Mutagenicit	Endo crine Disruption	Reproductive	Developmental	Repeat Exposure	Single Exposure	Systemic Kebeat Exposure	c Toxicity Single Exposure	Skin Sensitization	Skin Irritation	Eye Irritation	Acute Aquatic Toxicity	Chronic Aquatic Toxicity	Persistence	Bioaccumulation	Exposure
71-43-2 AIGBT Benzene	м	L	VH	VH	VH	н	н	н	н		н	н	I	н	н	н	м	М	н	н
75-01-4 AIGBT Vinyl chloride	М	L	I.	VH	VH	L	М	М	н	н	н		М	н	T	М	VH	М	L	VH
25265-71-8 GBT Dipropylene glycol	L	I.	L		L		L	L			L			М	н	L			L	
57-55-6 AIGBT 1,2-Propylene gly	L	- I	L	L	VH	н	L	Н	Н	Н	н		L	L	L	L	L	L	L	VH
Diethylene glycol	М	I	L	T	L	н	М	L		1	н	I.	L	L	L	L	L	L	L	VH
2-Butoxyethanol	м	м	м	L	VH	н	М	L		н	н	н	L	н	н	L	L	L	L	VH
115-11-7 GBTM Isobutene	1	L	I.	I.	L	L	I.	Н	L	T	н	I.	T	T	I.	L		М	L	VH
141-32-2     GBTM       Butyl acrylate	М	н	М	Т	L	L	М	L			М	М	н	н	н	н	н	L	L	Н
1322-13-0 Ethylhexyl acrylate	L				L	L		L								VH			I	
8002-74-2 GBT Hydrocarbon waxes	L	I.	L	T	L		L	I.	T		н	М	T	L	М	L	I.		L	
9003-07-0 GBT Polypropylene				I.												L			L	
9002-88-4 GBT Polyethylene AS I		М		Т												L			L	

### Data to Excel in <60s



AutoSave Off 🕞	5.6	- – – – – – – – Hazard Profilir	ng of East F	alestine che	micals.xlsx	🕜 No La	bel • Last I	Modified:	3m ago 🗸	م	Willia	ams, Antor	y (he/him/h	iis) 🌘	Γ	-	οx						
File Home In	sert Drav	v Page Layout Formula	s Data	Review	v Viev	v Dev	eloper	Help							🖓 Com	nents	🖻 Share						
Calibri Paste ~ ~ ~ Clipboard ~	U v Font	$\begin{array}{c c} 11 & A^{A} & A^{V} \\ \hline \\ 11 & A^{A} & A^{V} \\ \hline \\ 11 & A^{A} & A^{V} \\ \hline \\ \hline \\ \\ \\ \hline \\ \\ \hline \\ \\ \hline \\ \\ \hline \\ \\ \\ \hline \\ \hline \\ \\ \hline \\ \hline \\ \\ \hline \\ \\ \hline \\ \hline \\ \hline \\ \hline \\ \hline \\ \\ \hline \hline \\ \hline \\ \hline \\ \hline \hline \hline \hline \\ \hline \hline \hline \hline \\ \hline \hline \hline \\ \hline \hline \\ \hline \hline$	Y ~ a E → E Alignmer	b Wrap Text ∃ Merge & nt	: Center ~	General \$ ~ '	% 🤊 🗧		Conditiona Formatting	Format a • Table • Styles	s Cell Styles Y	E Inse Dele Form Cell	rt ~ ∑ te ~ ⊡ nat ~ ∢	Sort Filter	7 & Find & * Select *	Sensitiv Sensitiv	ity						
A1 • :	× ✓	fx The Hazard Compariso	n Das <mark>hbo</mark>	ard is a pr	ototype to	ool and a	compilati	on of inf	ormation s	ourced fr	om many	sites, dat	abases an	d sources	s including	g U.S. Fed	eral 👻						
А	В	С	D	E	F	G	Н	1	J	К	L	М	N	0	Р	Q	R 🔺						
The Hazard Comparis The data are not rev	son Dashboar iewed by USEI	d is a prototype tool and a comp PA – the user must apply judgme	ilation of i nt in use o	nformatior of the inform	sourced f nation. Th	rom many e results d	sites, data lo not indi	abases ar cate EPA'	d sources in s position o	ncluding U	S. Federal or regulati	and state on of thes	sources an e chemicals	d internat	ional bodi	es that sa	ves the us						
2         VH - Very High         H - High         M - Medium         L - Low           3           Human H											ffects	iciusive	NOL	dld	Authoritative								
3 4			Acute N	/ammalian	Toxicity		≩				Neurot	oxicity	Systemic Toxicity										
DTXSID	CAS	Name	Oral	Inhalation	Dermal	Carcinogenicity	Genotoxicity Mutagenici	Endocrine Disruption	Reproductive	Developmental	Repeat Exposure	Single Exposure	Repeat Exposure	Single Exposure	Skin Sensitization	Skin Irritation	Eye Irritation						
6 DTXSID3039242	71-43-2	Benzene	M		VH	VH	VH	н	н	н	н		н	н		н	н						
7 DTXSID8021434	75-01-4	Vinyl chloride	M	L	I	VH	VH	L	M	M	Н	Н	Н		M	Н	1						
8 DTXSID0027856	25265-71-8	Dipropylene glycol	L	I.	L		L		L	L			L			М	Н						
9 DTXSID0021206	57-55-6	1,2-Propylene glycol	L	I	L	L	VH	Н	L	Н	Н	Н	Н		L	L	L						
10 DTXSID8020462	111-46-6	Diethylene glycol	м	I	L	I.	L	Н	М	L		I	Н	I	L	L	L						
11 DTXSID1024097	111-76-2	2-Butoxyethanol	М	М	М	L	VH	Н	М	L		Н	Н	Н	L	Н	Н						
12 DTXSID9020748	115-11-7	Isobutene	1	L	1	1	L	L	I	Н	L	1	Н	I	I	1	I						
13 DTXSID6024676	141-32-2	Butyl acrylate	М	Н	М		L	L	М	L			М	М	Н	Н	Н						
14 DTXSID001031472	1322-13-0	Ethylhexyl acrylate	L				L	L		L							<b></b>						
Hazard P	rofiles Ha	zard Records +										1		• •									

### Perfect Example of FAIR Data and APIs



• We owe a lot to FAIR data and availability of information

• We curate a lot of our chemistry data using public resources such as PubChem, ChEBI, Common Chemistry and others

• The availability of Public APIs takes things to another level!

 We have been using the PubChem API to harvest data so we can build new applications, like the Safety Module

### Cheminformatics Safety Module Integrate multiple data streams...



Cheminformatics version: DEV, build: 202	<b>s Modu</b> 23-09-25 1	<b>les</b> 8:50:32 L	ЛС			<b>∲</b> +	IAZARD	ର SAF	ETY 🤇	Alert	s 🖬	PREDICT	1.0	PRE	DICT 2.0	🕄 S	EARCH	🕼 ST/	ANDARDIZ	ZE 🔞	TOXPR	INTS
↑ Search o	hemical	by Nam	e, CASRN	or DTXS	SID C	2							Show Structure Full						•	2 0	PDF	Ō
Chemical	Safety	Properties	Signal	Explosive	Flammable	Oxidizers	Compressed Gas	Corrosive	Acute Toxicity	Irritant	Health Hazard	Env. Hazard	NFPA 704	Fire Fighting	A ccidental Release Measures	Handling and Storage	stability and Reactivity	Transport Information	Regulatory Information	Other Safety Information	RQ Category	RQ In pounds (kilograms)
12125-02-9 GBTS Ammonium chloride	Ω	$\bigtriangleup$	Danger							<u>(!</u> )			200	i	i	i	i	i	i	i	D	5,000 (2,270)
1327-53-3 GB <sup>2</sup> Arsenic oxide (As2O3)	$\Box$	$\Delta$	Danger											i	i	i	i	i	i	i		
100-44-7 IGBTR Benzyl chloride	Ω	$\Delta$	Danger							()			321	i	i	i	i	i	i	i	В	100 (45.4)
94-75-7 AIGB 2,4-Dichlorophenox	$\bigcirc$	$\Delta$	Danger							<b>(!)</b>				i	i	i	i	i	i	i	В	100 (45.4)
25168-26-7 GB <sup>2</sup> 2,4-D isooctyl ester	$\bigcirc$	$\Delta$	Warning							()		<b>E</b>		i	i	i	i	i	i	i		
330-54-1 Diuron	$\sim$	$\Delta$	Warning							<b>(!)</b>		×2	100	i	i	i	i	i	i	i	В	100 (45.4)
IGBT 51-28-5 The Cheminformatics M information in one local	Aodules is	<b>R</b> a set of p ata are no	prototype moo ot reviewed by	Jules whic USEPA-	ch are using - the user m	a compilat ust apply ju	ion of inforr udgment in	nation sourc	ced from m formation.	hany sites, d The results	atabases do not inc	and sources	s including position on	U.S. Fed	eral and st	cate source on of these	es and inte	rnational be	odies that sa	Anves the use	^ er time b	y providing

### The CompTox API is now public https://api-ccte.epa.gov/docs/index.html



### **Computational Toxicology and Exposure Data APIs**

EPA's computational toxicology research efforts evaluate the potential health effects of thousands of chemicals. The process of evaluating potential health effects involves generating data that investigates the potential harm, or hazard of a chemical, the degree of exposure to chemicals as well as the unique chemical characteristics.

The APIs provided by EPA enable users to extract specific data from various databases and integrate them into their applications. These data are also available for download on our Data Download page

As part of EPA's commitment to share data, all of the computational toxicology data is publicly available for anyone to access and use. EPA's computational toxicology data is considered "open data", and thus all of the data are free of all copyright restrictions, and fully and freely available for both noncommercial and commercial use.

#### **Limited Access APIs**

An API Key is needed to access these APIs. Each user will need a specific key for each application. Please send an email to request an API key.

#### **Chemical APIs**



Access APIs for searching chemicals, files for chemical structures, and chemical details.





Access APIs for human and ecotoxicology data.

Documentation





Access APIs for chemical bioactivity data.

-ph

Documentation





- Modeling is essential to our research efforts we have models covering dozens of endpoints and continuing to expand
- Careful data assembly and curation is required
- OECD Modeling principles are guiding our future modeling
- CompTox Chemicals Dashboard provides public access to our various curated data streams, pre-predicted data and realtime prediction



- Modeling Team Dan Chang, Nate Charest, Charlie Lowe, Todd Martin
- Valery Tkachenko Cheminformatics Module
- Our DSSTox curation team
- Our SCDCD colleagues and DevOps team