# An Online Prediction Platform to Support the Environmental Sciences

Ann M. Richard\*, Chris Grulke, Kamel Mansouri, Richard Judson, Antony J. Williams

National Center for Computational Toxicology, U.S. Environmental Protection Agency, Research Triangle Park, NC

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

March, 13-17, 2016 ACS Spring Meeting, San Diego, CA

# EPA's National Center for Computational Toxicology (NCCT)



- ToxCast (EPA) & Tox21 (Multi-Agency)
  - screening >3800 (ToxCast) to >10K (Tox21) environmentally relevant chemicals across 10's to 100's of HTS assays
- ToxRef DB guideline animal toxicity study reference DB
- ExpoCast (DBs, models), Rapid Toxicity Assessments
  - rapid exposure estimates, integration of hazard and exposure (IVIVE) estimates for RA
- Public-facing, web-dashboards
  - facilitate access to & utility of EPA data with web tools

CHEMISTRY

Chemical databases

Chemical linkages Cheminformatics

Chemical structures

SAR/QSAR models

# NCCT (iCSS) Chemistry Dashboard

- Why build yet another public chemistry tool?
- Who is it being built for?
- How is it being built?
- What functionality & data will be included?

# Using Chemistry to improve linkages across EPA Resources



- NCCT ToxCast/EDSP Dashboards
- NCCT ExpoCast, CPCat, ACToR
- EPA Chemical Safety & Sustainability (CSS) research
- EPA Substance Registry System (SRS)
- EPA Integrated Risk Information System (IRIS)
- EPA Office of Pesticides (Inerts, Active Ingredients)
- EPA Endocrine Disruption Screening Program (EDSP)
- EPA Office of Water's Chemical Contaminants List
- EPA Office of Pollution Prevention & Toxics ChemView

Just lists (pdfs, tables, names, CAS, etc.)!!

### Chemistry Foundation to Support Multiple NCCT & EPA Projects





# **DSSTox Update**



#### DSSTox\_v1

LEARN THE ISSUES | SCIENCE & TECHNOLOGY | LAWS & REGULATIONS | ABOUT EPA

#### National Center for Computational Toxicology (NCCT)

Home
About DSSTox
Work in Progress
Frequent Questions
Structure Data Files
Central Field Definition
Apps, Tools & More
DSSTox Community
Site Map
Glossary of Terms
Help

You are here: EPA Home » Research & Development » CompTox » DSSTox

DSSTOX Distributed Structure-Searchable Toxicity (DSSTox) Database Network is a project of EPA's National Center for Computational Toxicology, helping to build a public data foundation for improved structure-activity and predictive toxicology

capabilities. The DSSTox website provides a public forum for publishing downloadable, structure-searchable, standardized chemical structure files associated with chemical inventories or toxicity data sets of environmental relevance. More



DSSTox Structure-Browser information Page

- manually curated 25K substance records
- EPA-focus, environmental tox datasets
- Emphasis on accurate CAS-namestructure annotations
- Public resource for high-quality structure-data files (SDF)

### DSSTox\_v2

- Convert DSSTox tables to MySQL
- Develop curation interface
- Implement cheminformatics workflow
- Expand chemical content
- Expand data inventories
- Web-services & Dashboard access





### EPA's Distributed Structure-Searchable Toxicity Database (DSSTox)





### Building the Chemistry Software Architecture



\$€PA

renommental Protection

# iCSS Chemistry Dashboard

Inited States Invironmental Protection

L



- Web interface supporting EPA's Chemical Safety & Sustainability research (iCSS)
- To be released April 2016, currently in beta-testing
- Will provide public access to DSSTox db
- UI in development only <4mo, building on substantial in-house data resources
- Initially name/ID-based searching only ...



**Chemistry Dashboard** 

Search a chemical by systematic name, synonym, CAS number, or InChIKey

Q

Single component search Ignore isotopes

Need more? Use advanced search.

# Componentizing API for Resolving Chemical Names

SEPA Environment<sup>2</sup> Protection

Q

Search a chemical by systematic name, synonym, CAS number, or InChIKey

- Look familiar?
  - ChemSpider, PubChem, ChEMBL, NCI Resolver ...
- ONE reusable EPA component API that searches for Names and CAS Registry Numbers...
- ... and SMILES, and InChIKeys, and recognizes invalid CAS Numbers, and converts SMILES when they aren't in DB, etc. etc.

# iCSS Chemistry Dashboard Releasing in April 2016

EPA United States Environmental Protection

EVELOPMEN

Sisphenol A	82	Formerly public su RDF/sen	DSSTox ( bstance IE	GSID, new u D, essential f	nique for
Searched by Synonym: Found 1	result for 'bisphenol A'.	TADI / SCII		applications	5
D 3D Q 🔟 🖻 📥 -	Intrinsic Pro	perties			
H <sub>3</sub> C CH <sub>3</sub>	Molecula Average Monoiso	r Formula: C15H16O2 Mass: 228.291 g/mol topic Mass: 228.11503 g/m	lol	Q Search in DSSTox	
		and the Company of the Company's statements			
Chemical External Lini	Citation	PubChem Articles	PubChem Patents	Comments	
Chemical Properties CSV Excel	Citation ks Synonyms PubChem Biological Ac	PubChem Articles	PubChem Patents	Comments	*
Chemical Properties CSV Excel Property Solubility	Citation  Ks Synonyms PubChem Biological Ac  Average (Exp.)  0.001 (1)	PubChem Articles stivities Range (Exp.)	PubChem Patents Average (Pred.)	Comments Range (Pred.)	-
Chemical Properties CSV Excel Property Solubility Melting Point	Citation iks Synonyms PubChem Biological Ac Average (Exp.) 0.001 (1) 154 929 (7)	PubChem Articles otivities Range (Exp.) 0.0005257 to 0.0005257	PubChem Patents Average (Pred.) 0.38 (2) 144 033 (3)	Comments Range (Pred.) 0.003675 to 0.7565	
Chemical Properties CSV Excel Property Solubility Melting Point Boiling Point	Citation Citation Citation Ks Synonyms PubChem Biological Ac Average (Exp.) 0.001 (1) 154.929 (7) 200.0 (1)	PubChem Articles otivities Range (Exp.) 0.0005257 to 0.0005257 153.0 to 158.0 200.0 to 200.0	PubChem Patents Average (Pred.) 0.38 (2) 144.033 (3) 348 95 (2)	Comments Range (Pred.) 0.003675 to 0.7565 131.8 to 158.0 334.4 to 363.5	*
Chemical Properties CSV Excel Property Solubility Melting Point Boiling Point	OH         Citation           iks         Synonyms         PubChem Biological Ac           Average (Exp.)         0.001 (1)           154.929 (7)         200.0 (1)           3.357 (3)         3.357 (3)	PubChem Articles etivities Range (Exp.) 0.0005257 to 0.0005257 153.0 to 158.0 200.0 to 200.0 3.32 to 3.431	PubChem Patents Average (Pred.) 0.38 (2) 144.033 (3) 348.95 (2) 3 524 (3)	Comments Range (Pred.) 0.003675 to 0.7565 131.8 to 158.0 334.4 to 363.5 3.205 to 3.727	*
Chemical Properties CSV Excel Property Solubility Melting Point LogP Atmospheric Hydroxylation Rate	OH         Citation           ks         Synonyms         PubChem Biological Ac           Average (Exp.)         0.001 (1)           154.929 (7)         200.0 (1)           3.357 (3)         N/A	PubChem Articles ctivities  PubChem Articles  Range (Exp.)  0.0005257 to 0.0005257  153.0 to 158.0  200.0 to 200.0  3.32 to 3.431  N/A	PubChem Patents Average (Pred.) 0.38 (2) 144.033 (3) 348.95 (2) 3.524 (3) 0.0 (1)	Comments Range (Pred.) 0.003675 to 0.7565 131.8 to 158.0 334.4 to 363.5 3.205 to 3.727 4.237e-11 to 4.237e-11	

SEPA

Search Chemical

Asency

to the strength a protection

Q.

### iCSS Chemistry Dashboard Releasing in April 2016

С



Chemical Pr	operties: M	elting Poir	it					•
				Ave	rage	Ran	ge	•
		Experi	mental	154	.929 (7)	153.	0 to 158.0	•
		Predic	ted	144	.033 (3)	131.	8 to 158.0	
CSV Excel								
Property	Raw Result	Mean Result	Minimu Result	m <sub>()</sub>	Maximu Result	m <sub>()</sub>	Result Unit	Result Type
Estimated MP (oC)	131.76	131.8	131.8		131.8		°C	predicte

- Original raw source result view
- Users can submit comments
- Source details provided
- Range of results possible for single chemical

Property	Result	Result	Result	Result	Unit	Туре	Source
Estimated MP (oC)	131.76	131.8	131.8	131.8	°C	predicted	EPI SUITE
Melting Point	153-158 ℃	155.5	153.0	158.0	°C	experimental	Alfa Aesar
Melting Point	154-157 ℃	155.5	154.0	157.0	°C	experimental	Merck Millipore
Melting Point	153-158 °C	155.5	153.0	158.0	°C	experimental	Alfa Aesar
Melting Point	155-158 °C	156.5	155.0	158.0	°C	predicted	J and K Scientifi
Melting Point	156 °C	156.0	156.0	156.0	°C	experimental	тсі
Meltina	153 °C	153.0	153.0	153.0	°C	experimental	Jean-Claude Br

### iCSS Chemistry Dashboard Releasing in April 2016



• Links external resources: EPA, NIH, property predictors



# Linkage to External Predictor, e.g. Chemicalize



#### Chemicalize.org (powered by ChemAxon)



# How are linkages created? e.g., CASRN



- ACTOR <u>http://actor.epa.gov/actor/GenericChemical?casrn=80-05-7</u>
- Toxcast Dashboard <a href="http://actor.epa.gov/dashboard2/#chemical/80-05-7">http://actor.epa.gov/dashboard2/#chemical/80-05-7</a>
- NIST Webbook <a href="http://webbook.nist.gov/cgi/cbook.cgi?ID=80-05-7">http://webbook.nist.gov/cgi/cbook.cgi?ID=80-05-7</a>
- TOXNET <u>http://toxnet.nlm.nih.gov/cgi-bin/sis/search2/r?dbs+toxline:@term+@rn+80-05-7</u>
- EPA SRS
   <u>http://iaspub.epa.gov/sor\_internet/registry/substreg/searchandretrieve/ad</u>
   <u>vancedsearch/externalSearch.do?p\_type=CASNO&p\_value=80-05-7</u>
- Some links pass InChIKeys for searching
- Some links pass SMILES for searching/prediction
- Some links pass a list of rank-ordered identifiers for searching (think Google Books, Pubmed, Google Scholar)

# Requires linked website be "chemically indexed"

# Creating Linkages to Data Sources ...



Creating linkages to indexed content is the easy part, ... creating accurate chemical-data listings is hard!

© FP

ommental Protection



- Challenges in assembling a database
  - Sourcing high quality data sorting wheat from chaff
  - How to mesh data together based on structure? On name? On CAS number? Other identifier?
  - Checking self-consistency of data?
  - Structure validation versus property value validation
     very different challenges

# Existing Online Property Prediction Platforms



- ACD/Interactive Laboratory (Ilab)
- Igor Tetko's OCHEM
- The OECD Toolbox
- etc...

Learn from others' examples, ... make open, transparent, reusable, ... tailor to EPA's needs!

# iCSS Chemistry Dashboard

- Deliver web-based platform hosting chemical property data of interest to "environmental toxicology & exposure" scientists (powered by DSSTox)
- Provide access to experimental data sets
   e.g. EPISuite's PHYSPROP: logP, BP, MP, WSol, etc.

### but...

"data quality" is a major concern & challenge a LOT of time was spent cleaning & curating data!



etall Protection

The EPI (Estimation Programs Interface) Suite™ is a Windows®-based suite of physical/chemical property and environmental fate estimation programs developed by EPA's and Syracuse Research Corp. (SRC).





#### Applicability Domain (AD) of original EPI Suite LogP Model (accuracy of the global model is an issue)



\* 280 cmpd cluster outside AD of EPI Suite -> Compared to new model predictions



- Collect "Open Data" for various endpoints of interest to EPA, from e.g.
  - PubChem
  - eChemPortal
  - Open Data Sources
- & Build QSAR prediction models
  - ✓ Using modern machine-learning approaches
  - $\checkmark$  Ensuring high quality data as inputs
  - ✓ Populating database with measured & predicted values
  - ✓ Making experimental data training sets available
  - ✓ Allowing real-time predictions (in the future)

# Open Data Set - example

200,000 melting points and 13,000 pyrolysis data points

Research article	Open Access
The development of models to predict melti point data associated with several hundred	ing and pyrolysis thousand
compounds mined from PATENTS	Journal of Cheminformatics
Igor V. Tetko $12^*$ , Daniel M. Lowe <sup>3</sup> and Antony J. Williams <sup>4</sup>	December 2016, 8:2

fig**share** 

- Modeling this much data is a challenge!
- Accessibility to data?
  - Figshare.com -
  - DataDryad.com
  - Institutional Repositories
  - Publishers?

L My data				
Create a new item	0 kB 💿	20 GB S	earch my data	
Melting Point & Pyrolysis Point Dataset	STATUS ~	TYPE 🗸	CREATED ↓	SIZE
Melting Point and Pyrolysis Point Data for Tens of Thousands of Chemicals		FILESET	9.12.2015 13:35	645.36 MB
Serving the medicinal chemistry community with Royal Society of Chemistry cheminformatics platforms	٠	PRESENTATIO	1.5.2015 23:46	10.47 MB
The Benefits of Participation in the Social Web of Science			1.5.2015 23:37	6.58 MB
Accessing Royal Society of Chemistry resources and making chemistry mobile			22.3.2014	13.46 MB

Antony Williams

# Integrating other EPA predictors



- ExpoCast
  - > near & far-field exposure models
- Environmental Fate Simulator (EFS)
   > air/soil/water distribution, biotransformation
- CERAPP
  - strogen receptor activity QSAR model
- T.E.S.T (in progress)
  - > phys-chem properties & toxicity endpoints

# e.g., T.E.S.T.



**\$EPA** 

now now mental Protection

	States Environmental Prot	ection Agency	Español	中文::	繁體版   中文: 简体版   Tiếng Việt   한국어
Learn the Issues	Science & Techno	logy Laws & Regulations	About EPA		Search EPA.gov
Related Topics:	Safer Chemical	s Research			Contact Us Share
Toxicity	Estimat	ion Software	Tool (TE	ST)	Option Fathead minnow LC50 (96 hr)
On this page: • QSAR Methodol • What's New in V • Prior Version Hi • System Require • Installation Inst • Publications • Get Email Alerts	logies Version 4.1? story ements cructions	<ul> <li>Downloadable desktop app</li> <li>Several phys-o endpoints pred</li> <li>Multiple QSAR methods emplo</li> <li>Multiple views</li> </ul>	Windows hem and to icted modeling oyed of data	x	Daphnia magna LC50 (48 hr) T. pyriformis IGC50 (48 hr) Oral rat LD50 Bioaccumulation factor Developmental Toxicity Mutagenicity Normal boiling point Vapor pressure at 25°C Melting point
The Toxicity Estimestimate the toxic (QSARs) methodo toxicity from the p molecular descript	ation Software Too ity of chemicals us logies. QSARs are hysical characteris ors). Simple QSAR	ol (TEST) was developed to a ing Quantitative Structure Ac mathematical models used to stics of the structure of chemi models calculate the toxicity	llow users to easily tivity Relationships predict measures cals (known as of chemicals using	of Ja	Density Surface tension at 25°C Thermal conductivity at 25°C Viscosity at 25°C Water solubility at 25°C Molecular Descriptors

# T.E.S.T. QSAR Model **Prediction Views**





Model # 1296

#### Model statistics

Parameter	Value
Endpoint	Fathead minnow LC <sub>50</sub> (96 hr)
r <sup>2</sup>	0.793
q <sup>2</sup>	0.733
#chemicals	101
Model	Model # 1296

#### Model fit results

11 10 9 8 Predicted toxicity 7 6 5 з 2 2 з 8 9 10 11 4 5 6 7 Experimental toxicity

Toxicity prediction results for 333-41-5 for Hierarchical clustering method

	Prediction results		
Endpoint	Experimental value CAS: 333-41-5 Source: <u>ECOTOX</u>	Predicted value <sup>a</sup>	Prediction interval
Fathead minnow LC <sub>50</sub> (96 hr) -Log(mol/L)	4.81	5.39	$4.54 \le Tox \le 6.24$
Fathead minnow LC <sub>50</sub> (96 hr) mg/L	4.70	1.23	$0.17 \le Tox \le 8.71$

<sup>a</sup>Note: the test chemical was present in the external test set.

#### Prediction results

Cluster model	Test chemical descriptor values	Prediction interval -Log(mol/L)	r²	q²	#chemicals
<u>1296</u>	Descriptors	6.010 ± 1.136	0.793	0.733	101
<u>1300</u>	Descriptors	5.458 ± 1.312	0.729	0.645	111
<u>1301</u>	Descriptors	5.136 ± 1.169	0.747	0.718	294
<u>1302</u>	Descriptors	4.922 ± 1.182	0.774	0.751	641

Cluster model predictions and statistics

	Cluste	ls with violate	d constraints	
Cluster Model	r <sup>2</sup>	q <sup>2</sup>	# chemicals	Message
<u>1121</u>	0.810	0.576	10	Rmax constraint not met
<u>1209</u>	0.799	0.574	11	Fragment constraint not met
1247	0.919	0.647	20	Fragment constraint not met
1264	0.869	0.781	22	Fragment constraint not met
1268	0.675	0.553	24	Fragment constraint not met

#### Descriptor values for test chemical

-	
-1	1
	-

#### A Platform for Predictions mmental Protection Developing a platform for hosting prediction models Coordinating efforts with other EPA prediction platforms and programs Exposure Assessment Tools and Models Search: All EPA This Area Recent Additions | Contact Us You are here: EPA Home \* Chemical Safety and Pollution Prevention \* Prevention, Pesticides & Toxic €EPA Estimation Program Interface (EPI) Suite oosure Home United States Environmental Protection Agency LEARN THE ISSUES | SCIENCE & TECHNOLOGY | LAWS & REGULATIONS | ABOUT EPA Risk Management Sustainable Technology You are here: EPA Home \* Research \* Risk Management Research \* Sustainable Technol übertool: web applications for ecological risk assessment (beta version) **Risk Management** » Quantitative Structure Activity Relationship Sustainable Technology **Quantitative Structure Activity Relationship** Introduction evaluating pesticide isks ecosystems.

LEARN THE ISSUES SCIENCE & TECHNOLOGY LAWS & REGULATIONS ABOUT EPA

#### Ecosystems Research

Methods, Models, Tool, &

**Ecosystems Research** 

Home

Databases

Publications

You are here: EPA Home » Exposure Research » Ecosystems Research » Environmenta

**Research in Action** 

Environmental Fate Simulator: Forecasting how chemicals move in the environment

# Conclusions

- SEPA Environment<sup>a</sup>l Protection
- iCSS Chemistry Dashboard will provide public access to data & services focused on chemicals of interest to EPA & environmental tox communities
- Initially serve up results for measured & prepredicted data, on-the-fly predictions in the future
- Data quality impacts models *public domain chemical-data linkages require curation/validation*
- All data and models will be available as OPEN DATA and OPEN CODE

# Stay tuned!!

# Acknowledgements



### NCCT Contributors to the iCSS Chemistry Dashboard

- Jeff Edwards
- Jeremy Fitzpatrick
- Jordan Foster
- Jason Harris
- Dave Lyons
- Kamel Mansouri
- Aria Smith
- Jennifer Smith
- Indira Thillainadarajah
- Chris Grulke
- Antony Williams



### Want to learn more?



1. Expansion of DSSTox: Leveraging public data to create a semantic cheminformatics resource with quality annotations for support of U.S. EPA applications (CINF 130)

Division of Chemical Information PAPER ID: 2385741

SESSION: Chemistry, Data & the Semantic Web: An Important Triple to Advance Science, 8:15 AM - 11:55 AM Wednesday, March, 16, 2016 from 10:40 AM - 11:05 AM Room 25B - San Diego Convention Center

2. Influence of data curation on QSAR Modeling – examining issues of quality versus quantity of data (MPPG 124)

Multidisciplinary Program Planning Group PAPER ID: 2394881

SESSION: Big Data Science, 1:30 PM - 4:30 PM Thursday, March, 17, 2016 from 1:30 PM - 2:00 PM ROOM & LOCATION: Room 3 - San Diego Convention Center