

DSSTox: The Open Environmental Chemistry Data underlying the CompTox Chemical Dashboard

<u>Chris Grulke</u>† Antony Williams Ann Richard

NCCT, U.S. EPA



OpenTox USA 2017

12-13 July 2017, Durham, NC

Office of Research and Development

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



NCCT's Purpose

- NCCT outputs: include a lot of data, models, algorithms and software applications
- We produce Open Data we want people to interrogate it, learn from it, develop understanding

Toxicity Forecasting

Advancing the Next Generation of Chemical Evaluation

EPA needs rapid and efficient methods to prioritize, screen and evaluate thousands of chemicals. EPA's Toxicity Forecaster (ToxCast) generates data and predictive models on thousands of chemicals of interest to the EPA. ToxCast uses high-throughput screening methods and computational toxicology approaches to rank and prioritize chemicals. In fact, EPA's Endocrine Disruption



Screening Program (EDSP) is working to use ToxCast to rank and prioritize chemicals.

- ToxCast has data on over 1,800 chemicals from a broad range of sources including industrial and consumer products, food additives, and potentially "green" chemicals that could be safer alternatives to existing chemicals.
- ToxCast screens chemicals in over 700 high-throughput assays that cover a range of high-

Downloadable Computational Toxicology Data

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.

As part of EPA's commitment to share data, all of the computational toxicology data is publicly available for anyone to access and use.

High-throughput Screening Data

EPA researchers use rapid chemical screening (called high-throughput screening assays) to limit the number of laboratory animal tests while quickly and efficiently testing thousands of chemicals for potential health effects.

• ToxCast Data: High-throughput screening data on thousands of chemicals.

Rapid Exposure and Dose Data

EPA researchers develop and use rapid exposure estimates to predict potential exposure for thousands of chemicals.

 <u>High-throughput toxicokinetics data</u>: It is important to link the external dose of a chemical to an internal blood or tissue concentration, this process is called toxicokinetics. EPA researchers measure the critical factors that determine the distribution.





Environmental Topics	Laws & Regulations	About EPA		Search EPA.gov	٩
TSCA Chemica	l Substance I	nventory	CONTACT U	JS SHARE f y (9

- Inventory was initially published in 1979
- Second version, containing about 62k chemical substances, was published in 1982
- Continues to grow and now lists ~85k chemicals, about 15k are confidential business information



Chemical representation levels supporting data integration





DSSTox History

GOAL: Linking chemical structures to data enabling SAR

- First release of data files in 2004
- Focused on high impact sets of data
 - -Carcinogenic Potency Database
 - -Drinking water disinfection by-products
 - -EPA's Integrated Risk Information System
 - -FDA's Maximum Daily Dose dataset
 - -EPA's Fat Head Minnow Toxicity dataset
 - -etc...
- Managed all chemical registration for ToxCast and Tox21 chemicals
- By 2014, roughly 20K manually curated substance records





DSSTox's Strengths and Weakness (circa 2014)

- Strengths
 - -Data as accurate as humanly possible
 - -Structure searching within subsets of data
 - -1:1:1 relationships between names, CASRNs, and structures
 - 3 table data model: list records, generic substances, compound structures
- Weaknesses
 - -Too many spreadsheets
 - –Too small
 - -Too manual



New Storage Architecture





DSSTox's Strengths and Weakness (circa 2014)

- Strengths
 - -Data as accurate as humanly possible
 - -Structure searching within subsets of data
 - -1:1:1 relationships between names, CASRNs, and structures
 - 3 table data model: list records, generic substances, compound structures
- Weaknesses
 - -Too many spreadsheets
 - –Too small
 - -Too manual



Example Data Load (EPASRS)

	EPA SRS												
Total Records		76944											
Internal Structure	All inchikeys agree			There is an agreement			There is only 1 inchikey			There is disagree ment	No in	chikey	
Спеск		16527		2815			11180			9702	36	620	
Other Source	agree	disagree	no structure	agree	disagree	no structure	agree	disagree	no structure		no smiles/nam e same	other	
Checks	15326	816	385	1923	719	173	5858	3598	1724		17378	19242	

SRS Internal Disagreement = 9702/76944 = 12.5% SRS External Disagreement =7415/30522 = 24.3%



Example Data Load (EPASRS)

Passing other			With S	Structure	9	No Structure				
source checks			23	3107			17378			
DSSTox Checks	agree	disagree	no structure data	inchi match different casrn	name mismatch	no record	same name (or synonym)	different name	name maps to different casrn	no record
	5703	304	46	144	8	16902	44	1109	72	16153

DSSTox / SRS Disagreement =502/6205 = 8.1%

WARNING: Always consider the accuracy/consistency of chemical data

EPA United States Environmental Protection Agency



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



Keeping Manual Curation Alive

View/Edit a Structure Sea Single Record	rch Browse/Curate Export DSSTox Chemotypes Records	Manage Manage Property Add [Chemical Lists Data Casrn	Deleted ns	
Preferred Name matched null You are viewing the record associated with DTXSID80198757 CASRN: 62885-41-0 4-Hydroxy-3-methox	□ □ □ ○ × □ ① ●	N → CH ³	H C N O S F P C I Br I I	
	Calculate from Structure Substance_ID: DTXSID80198757 CAS: 62885-41-0 Name: 4-Hydroxy-3-methoxypyridine Substance Type: Single Compound ▼ QC Level: DSSTox_High ▼ Data Source: STN(DSSTox) ▼ CAS [50700-60-2] assigned by DSSTox to pyridin-one tautom form, which resolves to hydro form thru InChI	Compound_IC Chemical Sho Private Notes Source of CAS Double Stereo Chiral Stereo Chemical Form Organic Form	D: DTXCID40121248 Dwn: Tested Chemical S: S-Compound: STN(DSSTox) ▼ o: None ▼ : None ▼ m: Organic ▼	• • • •



W E

Internal List Conflicts: ChemReg List Curation

ew/Edit a Structure agle Record	Search Bro Rec	wse/Cui ords		Export DSSTox	Chemotypes	Manage Chemical Lists	Manage Property Data	Add D Casrn	Deleted Is	
/elcome cgrulke		L	.ist:	ALANWOOD						
diting Listname: ALANW	/OOD						(1 of 69)		1 2 3 4 5 6 7 8 9 10 🕨 🖬 25	
Internal Check R	esults			Record ID		Externa	l ID		1st Identifier	Warning
Description	Records		0	DTXRID303936283	3 (3-etho	xypropyl)mercury	bromide		6012-84-6	NONE
All	1718									
Duplicate CASRN	2								Identifier	
Duplicate			(3-	-ethoxypropyl)mercur	y bromide					
STRUCTURE, Duplicate STRUCTURE INCHIKEY	2		60	12-84-6						
Invalid casrn	2		Int	ChI=1S/C5H110.BrH.	.Hg/c1-3-5-6-	4-2;;/h1,3-5H2,2H3	3;1H;/q;;+1/p-1			
NONE	1712		BW	UIOGHGUVLNSX-UH	FFFAOYSA-M					
		' -		DTVDID000002000	a <u>a a</u> 4 a	L			70.07.5	NONE
		_	•	DTXRID003936284	+ 1,2-dici	nioropropane			/8-8/-5	NONE
			0	DTXRID703936285	5 1,3-dic	hloropropene			542-75-6	NONE
			0	DTXRID403936286	5 1-meth	ylcyclopropene			3100-04-7	NONE
			0	DTXRID103936287	7 1-naph	thol			90-15-3	NONE
			0	DTXRID803936288	8 2-(octy	lthio)ethanol			3547-33-9	NONE
			0	DTXRID503936289	9 2,3,5-ti	ri-iodobenzoic acid			88-82-4	NONE
			0	DTXRID803936290	0 2,3,6-T	ВА			50-31-7	NONE
			0	DTXRID503936291	1 2,4,5-T				93-76-5	NONE
			0	DTXRID203936292	2 2,4,5-T	В			93-80-1	NONE
			0	DTXRID903936293	3 2,4-D				94-75-7	NONE



ChemReg List Curation (cont.)

View/Edit a Structure Single Record	e Search Br Re	rowse/C ecords	Curate	Export DSSTox Chemotypes	Manage Manage Pro Chemical Lists Data	perty Add Deleted Casms		Wel	come, Chris
Welcome cgrulke						Substance Ma	apping		
Editing Listname: ALAN	VOOD				(1	of 5) 📧 🔍 📘 2 3	4 5 🕨 🕨 25 🔻		
External Check F	Results			Source Casrn	Source Name	Hit Substance_ID	Hit Casrn	Hit Name	
Description	Records		0	88-82-4	2,3,5-tri-iodobenzoic acid	DTXSID4041317	88-82-4	2,3,5-Triiodobenzoic acid	Validate Ma
Resolved Duplicates	0		0	50-31-7	2,3,6-TBA	DTXSID6040296	50-31-7	2,3,6-Trichlorobenzoic acid	Validate Ma
Ignored Structure matched	0		0	122-88-3	4-CPA	DTXSID9034282	122-88-3	4-Chlorophenoxyacetic acid	Validate Ma
Preferred Name matched NAME	2		0	126448-41-7	acibenzolar	DTXSID20155187	126448-41-7	Acibenzolar [ISO]	Validate Ma
CAS-RN matched CASRN			0	76636-10-7	amibuzin	DTXSID20227459	76636-10-7	Amibuzin [ISO]	Validate Ma
Structure matched STRUCTURE			0	3566-10-7	amobam	DTXSID0058067	3566-10-7	Ambam	Validate Ma
matched NAME CAS-RN matched	71		0	86-88-4	antu	DTXSID8020919	86-88-4	1-(1-Naphthyl)-2-thiourea	Validate Ma
CASRN Structure matched STRUCTURE Unique Synonym matched NAME	106		0	52-46-0	apholate	DTXSID7073149	52-46-0	1,3,5,2,4,6- Triazatriphosphorine, 2,2,4,4,6,6-hexakis(1- aziridinyl)-2,2,4,4,6,6- hexahydro-	Validate Ma
CAS-RN matched CASRN			0	3586-60-5	asomate	DTXSID70189412	3586-60-5	Arsine, tris(dimethyldithiocarbamoy	Validate Ma
Structure matched STRUCTURE			0	28956-64-1	bentaluron	DTXSID30183153	28956-64-1	Bentaluron [ISO]	Validate Ma
matched NAME Other CAS-RN	2		0	21564-17-0	benthiazole	DTXSID6032647	21564-17-0	2- (Thiocyanomethylthio)benzo	Validate Ma
matched CASRN Structure matched			0	1022-46-4	bentranil	DTXSID60144732	1022-46-4	4H-3,1-Benzoxazin-4-one, 2-phenyl-	Validate Ma



Why Curation is Essential





Chemical Families (e.g. PCBs)

Successor Substances (209)





Next Steps: How to Work with UCVBs

Xylene

▼ S	ucces	sor Substances (3)			
		CAS-RN	Relationship	Source	St
	•	108-38-3	is a Representative Isomer of this	STN(DSSTox) ▼	~
	•	95-47-6	is a Representative Component o	STN(DSSTox) V	~
		106-42-3	is a Representative Component o	STN(DSSTox) V	~
			Delete Selecte	ed Add Related Cas	

Alcohols, C6-12, ethoxylated









• 3-Layer data model



- Be careful with public chemical data
- Please take and use our data!



Acknowledgements



EPA NCCT Imran Shah Chris Grulke Jeff Edwards Ann Richard Jordan Foster Jennifer Smith Richard Judson Grace Patlewicz John Wambaugh Michelle Krzyzanowski