Advancing Non-Targeted Analysis of Xenobiotics in Environmental & Biological Media



EPA's DSSTox Chemical Database: A Resource for the Non-Targeted Testing Community

Ann Richard National Center for Computational Toxicology Office of Research & Development, US EPA

> August 18-19, 2015 Research Triangle Park, NC

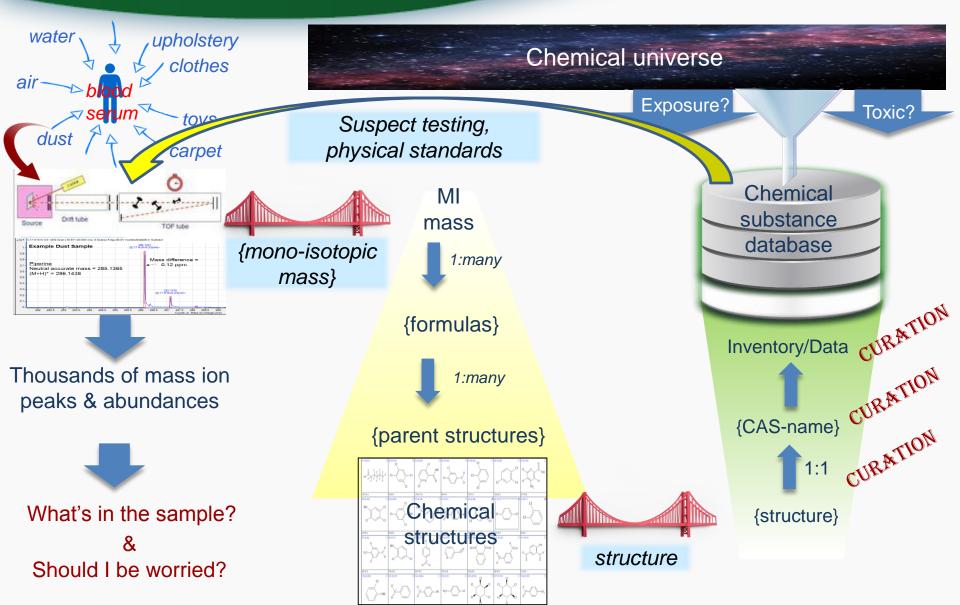
Outline



- Cheminformatics view of problem
- DSSTox chemical database
- ToxCast chemical library
- □ Tox21 analytical QC
- □ Challenges

Cheminformatics view of non-targeted testing problem





How big is the problem?



http://www.chemspider.com/blog/

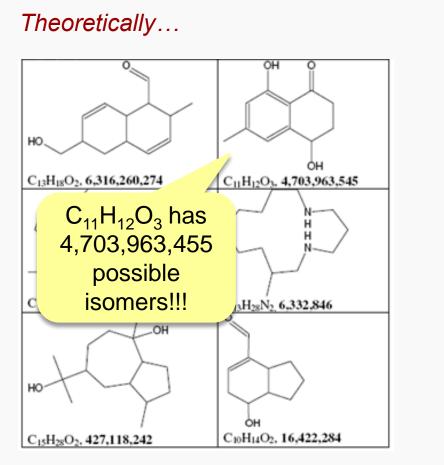


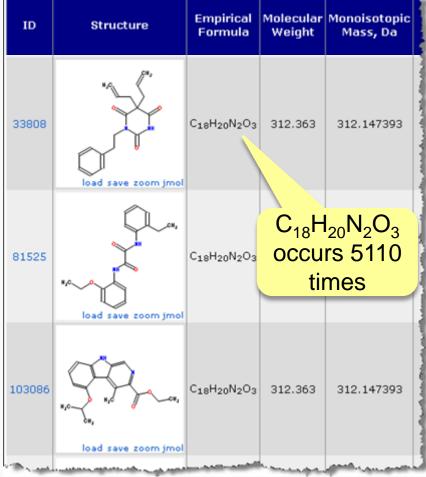
How Many Structures Can You Generate From A Molecular Formula?

Posted by: Antony Williams in Chem Spider Chemistry (22 million ChemSpider IDs in 2008)

Copyright©2008 Antony Williams

Highest frequency formula ChemSpider?





How big is the problem?



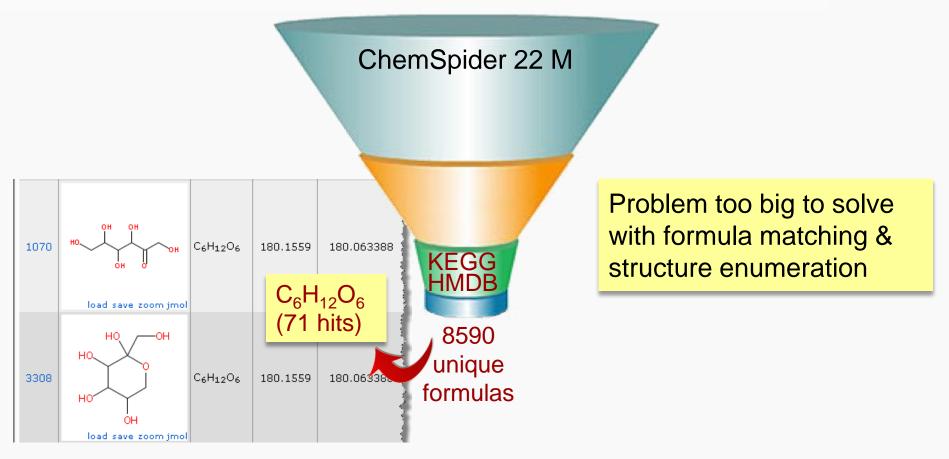
http://www.chemspider.com/blog/



How Many Structures Can You Generate From A Molecular Formula?

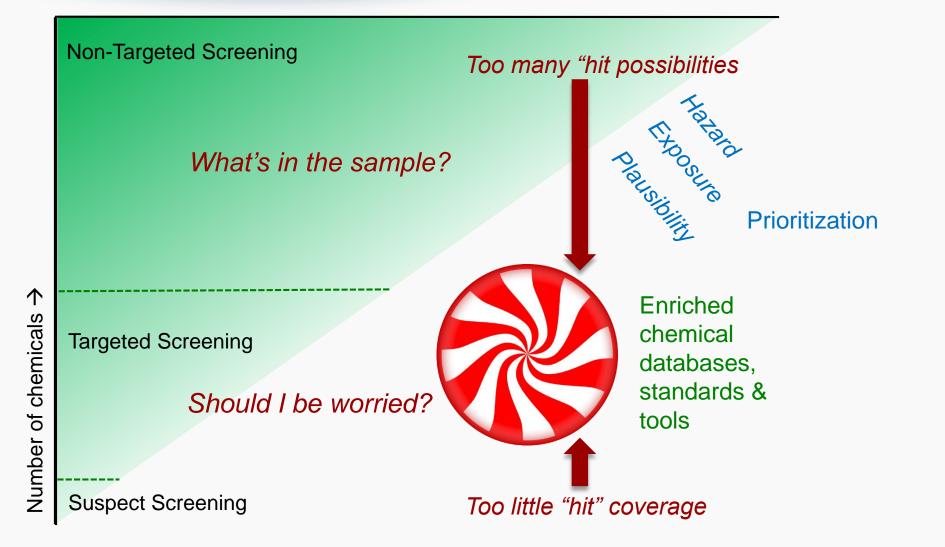
Posted by: Antony Williams in Chem Spider Chemistry (22 million ChemSpider IDs in 2008)

Copyright©2008 Antony Williams



Cheminformatics view of non-targeted testing problem





DSSTox_v1 (thru 3/2014)



€ FP4 ed States Environmental Protection Agency LEARN THE ISSUES | SCIENCE & TECHNOLOGY | LAWS & REGULATIONS | ABOUT EPA National Center for Computational Toxicology (NCCT) You are here: EPA Home » Research & Development » CompTox » DSSTox Home DSSTox About DSSTox Work in Progress Distributed Structure-Searchable Toxicity (DSSTox) Database Network is a project of EPA's National Center for Computational Toxicology, helping to buil Frequent Questions public data foundation for improved structure-activity and predictive toxicolog Structure Data Files capabilities. The DSSTox website provides a public forum for publishing downloadable, structure-searchable, standardized chemical structure files **Central Field Definition Tabl** associated with chemical inventories or toxicity data sets of environmental Apps, Tools & More relevance. More DSSTox Community Site Map Structure-Browser Glossary of Terms Help DSSTox Publications DSSTox Structure-Browser information Page

http://www.epa.gov/ncct/dsstox/

- Original target audience: Structure-Activity Relationship (SAR) toxicity modeling community
- Focus on EPA, HPV, environmental toxicity datasets
- Emphasis on accurate CAS-name-structure annotations at substance level
- Public resource for high-quality structure-data files
- ARYEXP_v2a_958_06Mar2009 CPDBAS_v5d_1547_20Nov2008 DBPCAN_v4b_209_15Feb2008 EPAFHM_v4b_617_15Feb2008 FDAMDD_v3b_1216_15Feb2008 GEOGSE_v2a_1179_09Mar2009 HPVCSI_v2c_3548_15Feb2008 HPVISD_v1b_1006_15Feb2008 IRISTR_v1b_544_15Feb2008 KIERBL_v1a_278_17Feb2009 NCTRER_v4b_232_15Feb2008 NTPBSI_v4c_2330_04Aug2009 NTPHTS_v2c_1408_11Mar2009 TOX215_v2a_8193_22Mar2012 TOXCST_v4a_1892_20Mar2012* External Data Files (ISSCAN)

DSSTox_RID
DSSTox_GSID
DSSTox_CID
DSSTox_FileID
TestSubstance_ChemicalName
TestSubstance_CASRN
TestSubstance_Description
ChemicalNote
STRUCTURE_Shown
STRUCTURE_Formula
STRUCTURE_MolecularWeight
STRUCTURE_ChemicalType
STRUCTURE_TestedForm_DefinedOr
ganic
STRUCTURE_ChemicalName_IUPAC
STRUCTURE_SMILES
STRUCTURE_Parent_SMILES
STRUCTURE_InChIS
STRUCTURE_InChIKey
Substance_modify_yyyymmdd

Approx. 25K CAS-substances, 16K structures

DSSTox Update



DSSTox_v1

SEPA United States Environmental Protection Agency

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

National Center for Computational Toxicology (NCCT)

Home

Structure Data Files

Apps, Tools & More

DSSTox Community

Glossary of Terms

Site Map

Central Field Definition Table

About DSSTox Work in Progress Frequent Questions DSSTox Distributed Structure-Searchable Toxicity (DSSTox) Database Network is a

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

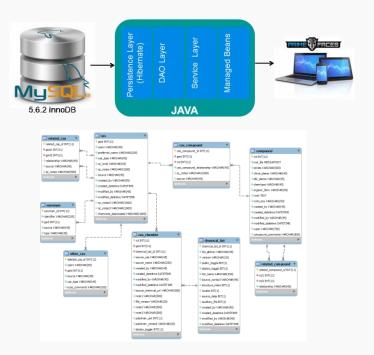
project of EPA's National Center for Computational Toxicology, helping to build a public data foundation for improved structure-activity and predictive toxicology capabilities. The DSSTox website provides a public forum for publishing downloadable, structure-searchable, standardized chemical structure files associated with chemical inventories or toxicity data sets of environmental relevance. More



- Website to be retired 9/30/2015
- DSSTox_v1 files & documentation will remain available on EPA ftp site
- Original 25K substance records stored in MS ACCESS and Excel tables serve as the initial Level 1,2 input to DSSTox_v2

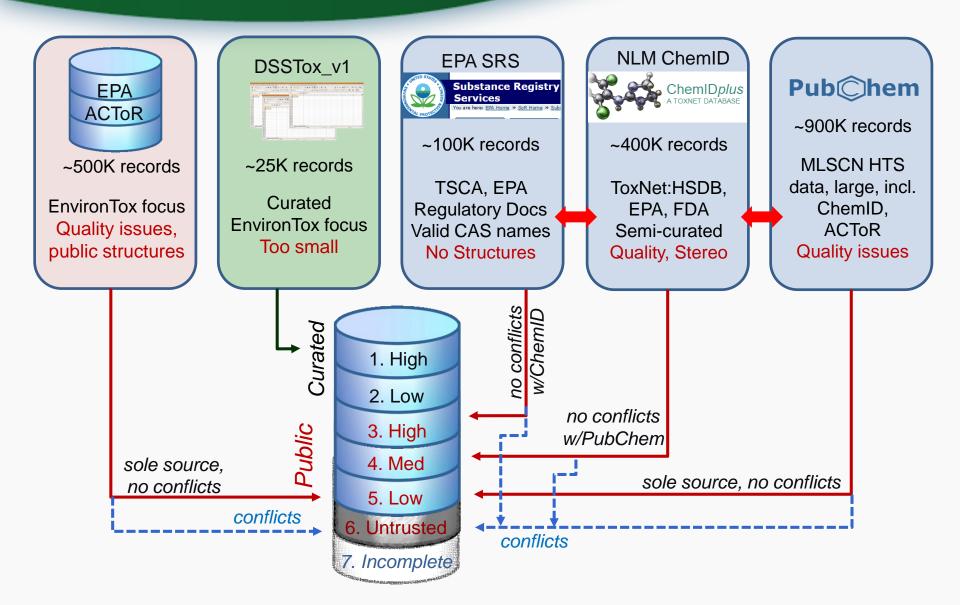
DSSTox_v2

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



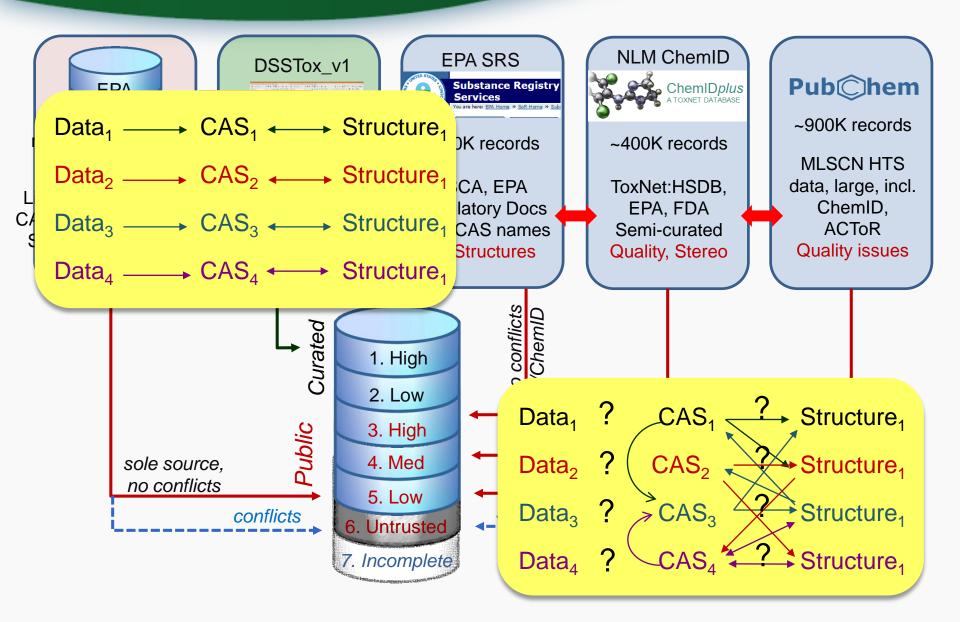
DSSTox_v2 CAS-Structure Sources: QC levels





DSSTox_v2 CAS-Structure Sources: QC levels





DSSTox_v2 Construction



Data source load order:

- 1) DSSTox_v1 (~22K)
 - ✓ 1:1 CAS-structure mappings
 - ✓ Assign NOCAS_GSID
 - ✓ Related CAS & structure mappings (e.g., NOCAS, mixtures)

2) EPA SRS (~77K)

- \checkmark systematic name \rightarrow structure conversion
- ✓ internal CAS-structure conflicts (12.5%)
- ✓ ChemID conflicts (24% of 30K overlaps)
- ✓ DSSTox conflicts (8% of 6200 overlaps) → queue for curation

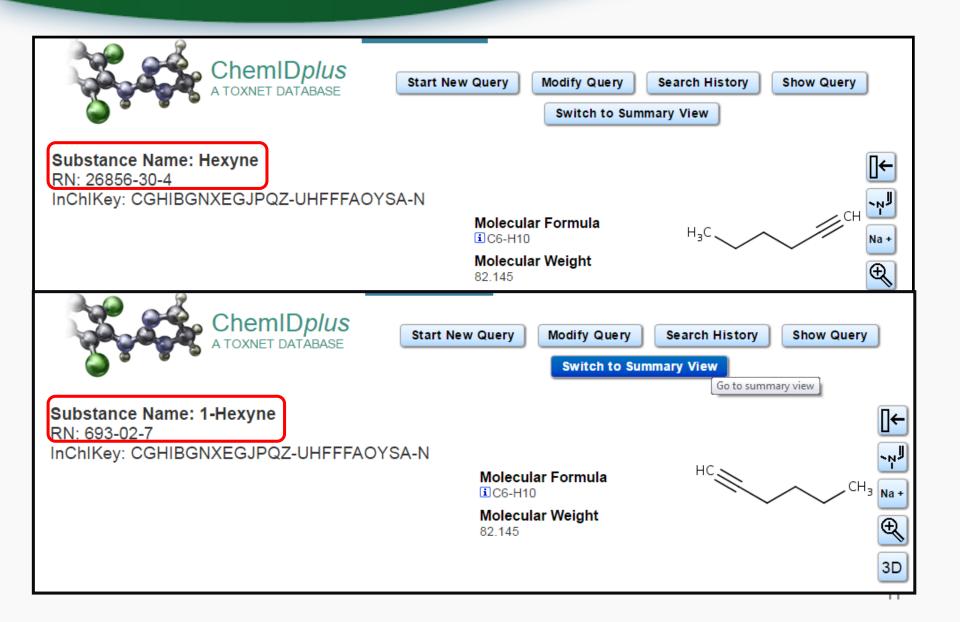
3) ChemID (~77K)

- ✓ internal CAS-structure conflicts (4.5%)
- ✓ PubChem conflicts (45% of 225K overlaps) ... OUCH!!
- ✓ DSSTox conflicts (11% of 2300 overlaps) → queue for curation

4) And so on ...

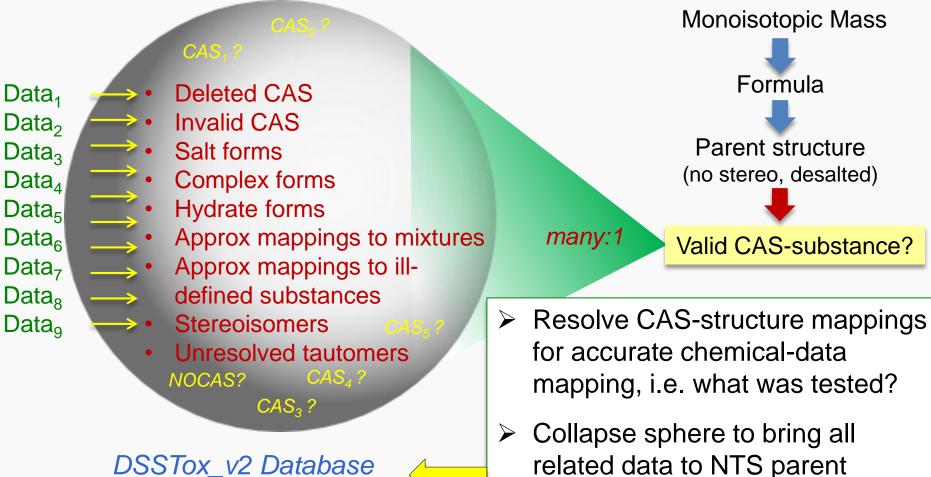
Example problem





CAS-Structure "Sphere of Confusion"





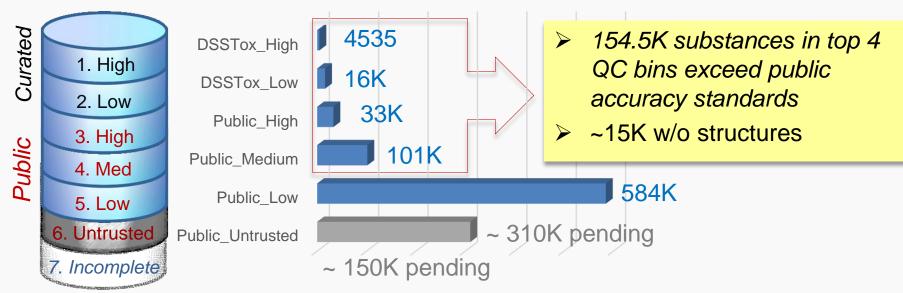
structure-formula level

& Cheminformatics Layer

DSSTox_v2 Totals



QC Level Totals (12Jun2015)



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

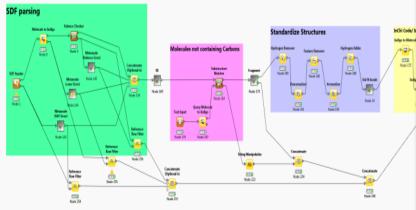
KNIME structure-"cleaning" workflow



https://www.knime.org/knime

Objectives:

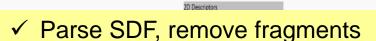
- Combine community approaches to structure processing
- Develop a flexible workflow to be used by EPA and shared publicly
- Process DSSTox files to create "QSAR-ready" structures



Publicly available cheminformatics toolkits in KNIME:







- ✓ Explicit hydrogen removed
- ✓ Dearomatization
- Removal of chirality info, isotopes and pseudo-atoms
- Aromatization + add explicit hydrogen atoms
- ✓ Standardize Nitro groups
- ✓ Other tautomerize/mesomerization
- ✓ Neutralize (when possible)



Slide courtesy of K. Mansouri

DSSTox Viewer (EPA Intranet)



DSSTox Viewer	http://rtpavaki1.epa.gov:8080/DSSToxViewer/						
Search a Single Record							
Welcome Matched Name. You are viewing the record associated with GSID: 20182 CASRN: 80-05-7	SynonymSourceBisphenol A (BPA) (4,4'-Propane-2,2- diyldiphenol) (Phenol, 4,4'-(1- methylethylidene)bis-)DSSTox						
GSID or CAS or	Other Cas (5) CAS-RN Relationship						
Name bisphenol A search	137885-53-1 Deleted SRS 146479-75-6 Deleted SRS						
GSID: 20182 CID: CAS: 80-05-7 Chemica Name: Bisphenol A Shown:	27360-89-0 Deleted SRS 28106-82-3 Deleted SRS						
QC Notes:	37808-08-5 → N → Since → Si						
QC Levels	al Form: Organic						
DSSTox_High: Hand curated and validated DSSTox_Low: Hand curated and confirmed using m Public_High: Extracted from EPA SRS and confirm Public_Medium: Extracted from ChemID and confirm Public_Low: Extracted from ACToR or PubChem Public_Untrusted: Postulated, but found to have conflic	ned to have no conflicts in ChemID and PubChem ned to have no conflicts in PubChem						

DSSTox Batch Tool (EPA Intranet)

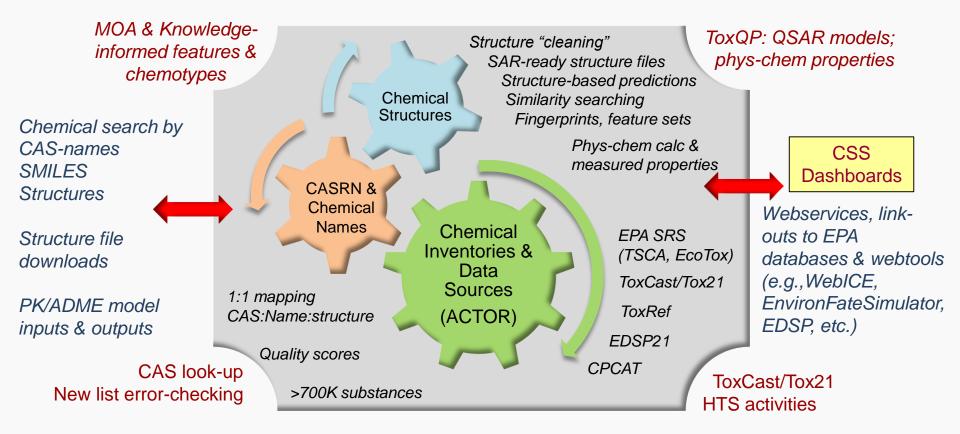


http://rtpavaki1.epa.gov:8080/DSSToxViewer/

DSSTox Viewe	er	ľ	Calibri	- 11 - A A	≡≡≡	- %	F Wrap Te	ext	General		-	₽ ₽
Search a Export Single Record DSSTox	Par	aste v ♥ B I U v	• 🖾 • <u> •</u> • <u>A</u> •		€ €	🔛 Merge	& Center 🔹	- \$ - %	ć ୭ €.0 .00	→.0	ditional Forr atting ▼ Ta	
	Clir	ipboard 🕞	Font 🕞		Alignn	ment	r	ra Nu	umber	G I	Style	
Choose your identifier t Name Submit a list of Identifie		B2	2 • : 🗙	√ f_x Match	ched Name.							
Ecolyst Valifenalate	^		А	В	с	D	E	F	G	н	1	L
Empenthrin [(EZ)- (16	R)-	1	Query	Found_By	DSSTox_G	DSSTox_C	ChemNan	CASRN	QC_Leve	I Descript	tic ChemNo	ot STRUCTU
isomers]	2/			Matched Name.	32628	12628	8 Methoxyf	161050-58	· -	· ·		C22H28N
Pyrafluprole Fluoxastrobin		3	Fenamidone	Matched Name.	34590	14590) Fenamido	161326-34	4 DSSTox_I	Single C	ompound	C17H17N
Pyriprole		4	(S)-Dimethenamid	Not Found								
Benthiavalicarb		5	Selamectin	Matched Name.	45903	25903	Selamecti	165108-07	7 DSSTox_I	Single C	or tautomer	ers C43H63N
Pyribencarb Triclopyricarb		6	Dinotefuran	Matched Name.	34549	14549	Dinotefur	165252-70	J DSSTox_I	Single C	ompound	C7H14N4
Meperfluthrin		7	Lepimectin	Matched Name.	58236	32043	Lepimecti	171249-05	5 DSSTox_F	A Single C	or Absolute	e C41H53N
Pyrametostrobin		8	Indoxacarb	Matched Name.	32690	12690	Indoxacar	173584-44	4 DSSTox_F	A Single C	ompound	C22H17C
Heptafluthrin Zeta-Cypermethrin	-	9	Pyraclostrobin	Matched Name.	32638	12638	8 Pyraclostr	175013-18	8 DSSTox_I	Single C	ompound	C19H18C
Zeta- <u>Cypermethrin</u>	10	Ethiprole	Matched Name.	58003	31771	L Ethiprole	181587-01	1 DSSTox_F	A Single C⁄	ompound	C13H9Cl	
		11	Acetoprole	Matched Name.	58017	31785	6 Acetoprol	209861-58	3 DSSTox_F	H Single C/	ompound	C13H10C
Download DSSTox	Gsids	12	Clothianidin	Matched Name.	34465	14465	Clothianic	210880-92	2 DSSTox_I	Single C	or tautomer	rs C6H8CIN
Download DSSTox Structures	13	Profluthrin	Matched Name.	58023	31791	l Profluthri	223419-20	J DSSTox_I	H Single C	ompound	C17H18F	
	14	Cloflubicyne	Matched Name.	274166	196562	2 Cloflubicy	224790-70	J DSSTox_I	Single C	or	C11H6Cl	
	15	Metofluthrin	Matched Name.	34738	14738	8 Metofluth	240494-70	J DSSTox_I	Single C	ompound	C18H20F	
	16	Pinoxaden	Matched Name.	34823	14823	Pinoxader	243973-20	J DSSTox_I	Single C	ompound	C23H32N	
	17	Dimefluthrin	Matched Name.	58043	31811	L Dimefluth	271241-1/	4 DSSTox_F	A Single C	ompound	C19H22F	
		18	Ecolyst	Not Found								
		19	Valifenalate	Matched Name.	58051	31819	Valifenala	283159-90	J DSSTox_I	A Single C⁄	or Absolute	a : C19H27C
		20	Empenthrin [(EZ)- (1R)	Matched Name.	58243		Empenthr	303021-82	2 DSSTox_F	H Mixture	/F Absolute	e stereoch
		21	Pyrafluprole	Not Found								
1		4 '	1		·	/	1					

CSS Chemical Explorer Dashboard (powered by DSSTox & ACToR)

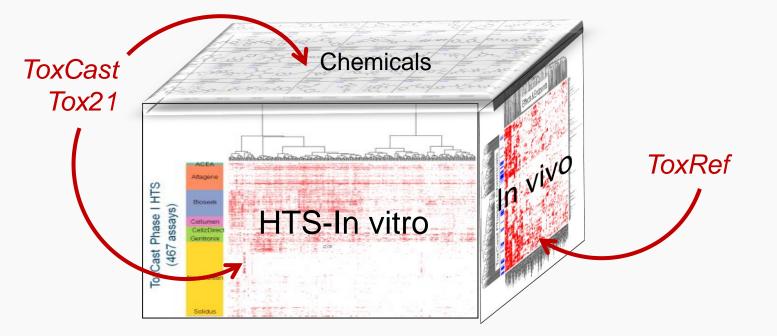




EPA's ToxCast/Tox21 Projects



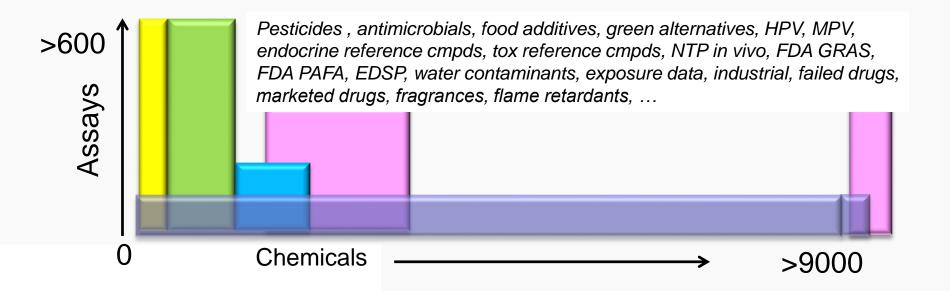
- Build a diverse, highly prioritized chemical library of interest to EPA regulatory programs (e.g., EDSP) and of relevance to environmental toxicology
- Use high-throughput screening (HTS) to generate bioassay profiles and fill data gaps for thousands of chemicals
- Use all of these data to improve ability to model adverse outcomes



ToxCast & Tox21 Inventories

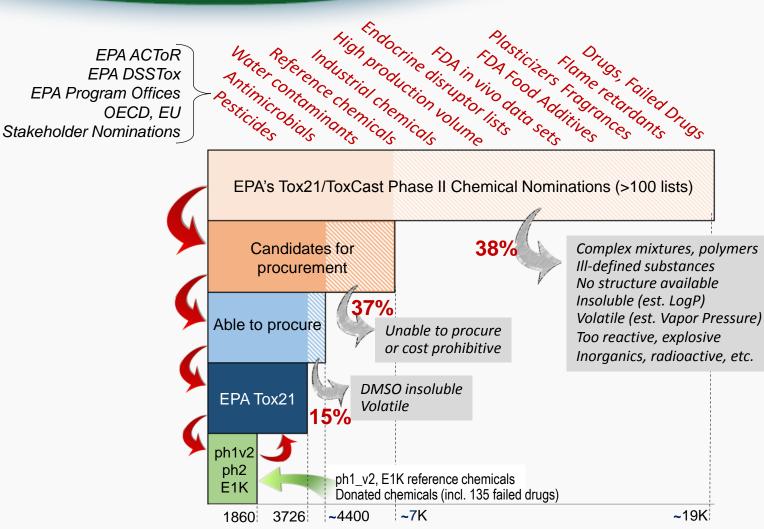


Set	Chemicals	Assays	Endpoints	Completion	Available
ToxCast Phase I	293	~600	~700	2011	Now
ToxCast Phase II	767	~600	~700	03/2013	Now
ToxCast E1K	800	~50	~120	03/2013	Now
Tox21	~8900	~80	~150	Ongoing	Ongoing
ToxCast Phase III	~2000	~300	~300	In process	2014-2015



Construction of EPA's Tox21/ToxCast Inventories

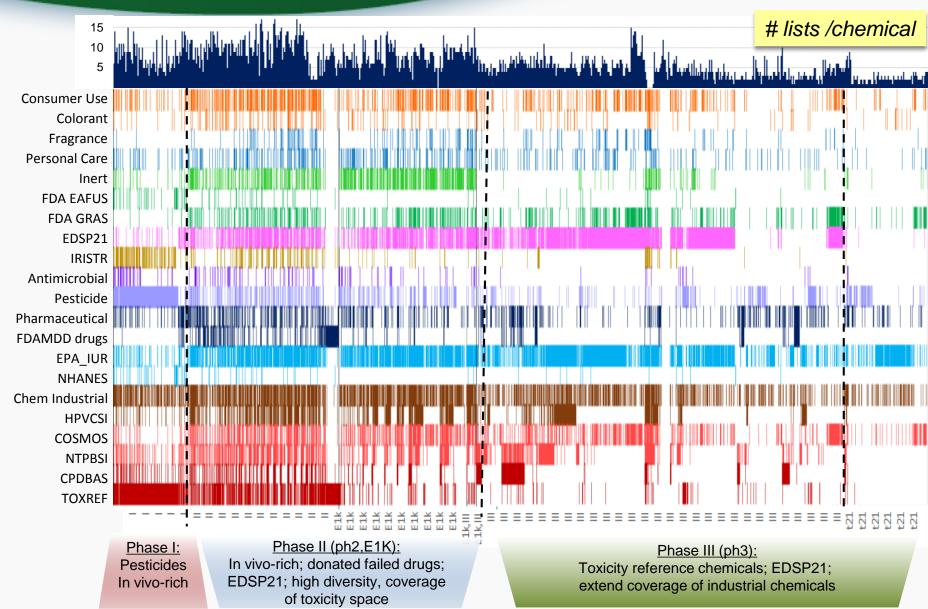




NUMBER OF CHEMICALS

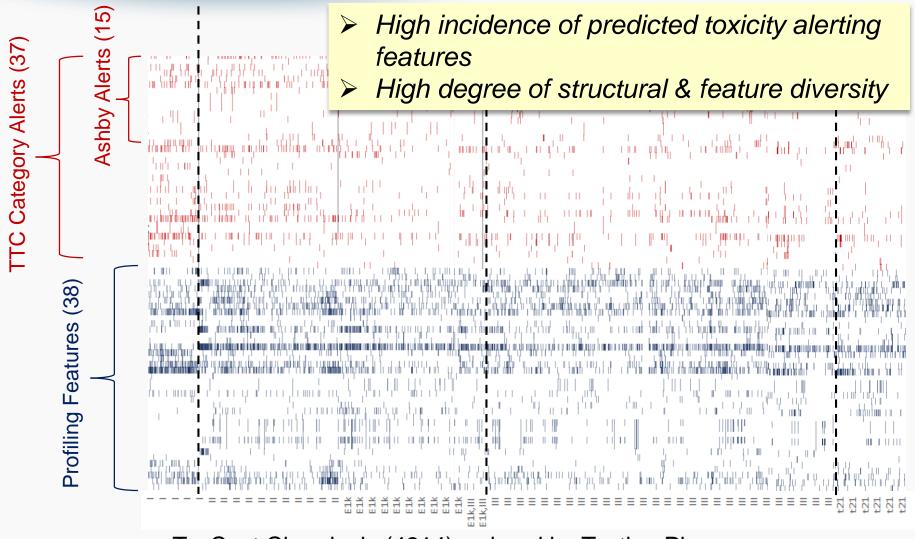
ToxCast Chemical Coverage: Use, Exposure, Toxicity





ToxCast: Toxicity Structure-Alerts & Generic Feature Coverage

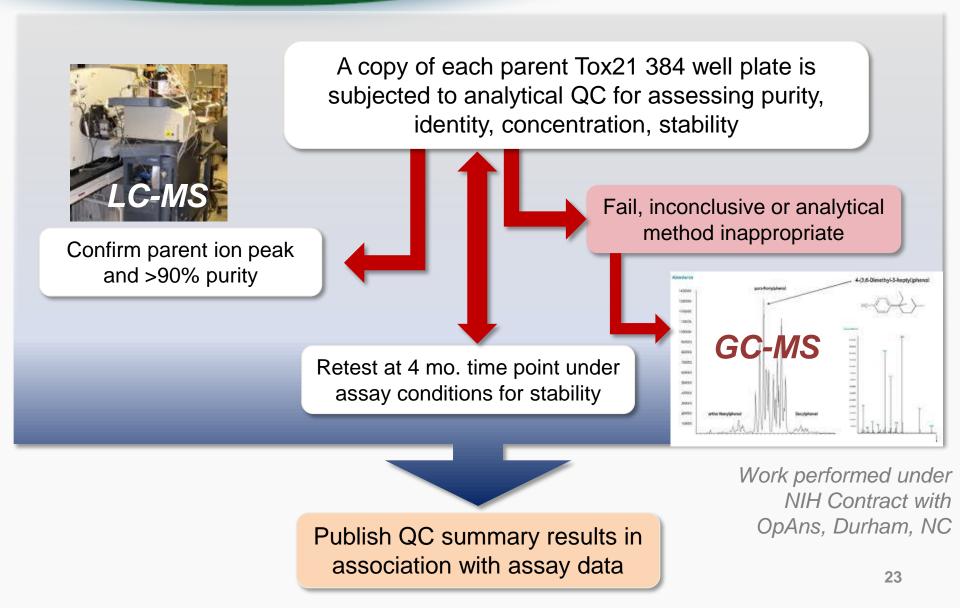




ToxCast Chemicals (4214) ordered by Testing Phase

Tox21 Analytical QC

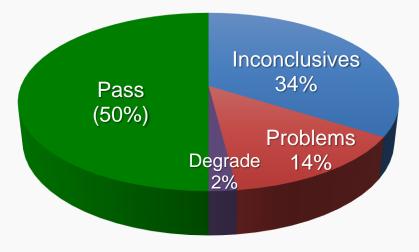




Tox21 and ToxCast Chemical Library Analytical QC Results (8/2015)



Tox21_QC_Sum-GSID (8593 total)



- 50% pass purity/ID/concentration checks
- A third(34%) of library pose analytical QC challenges (LC-MS and GC-MS)
- 2% degrade after 4 months under testing conditions
- 14% problems purity (<75%), ID and/or low concentration (<30% of expected [C])

> Which chemicals have QC issues? (e.g. SVOCs?)

- > Which chemicals were not analyzed? (e.g., mixtures, inorganics, etc.)
- ➤ How are the HTS activity profiles linked to QC?

Tox21 Analytical Chemical QC: Publicly available in PubChem



nmary.cgi?sid=144206248	
	х
	_
C L https://tripod.nih.gov/tox21/QC/Tox21_112695	
ID Tox21_112695 Plate FDA-plate8 Well P2-L-06 File00181048-01.D Inj Date: 15 Feb 10 1:47 pm • MF C19H20N2O3.CH4O3S.H2O MW 324.1 Expected Conc: 3.00 m	
*ELSD Norm. * <	
AMS AMS <th></th>	
1 1	
	 https://tipod.nih.gov/to: X https://tipod.nih.gov/tox21/QC/Tox21_112695 https://tipod.nih.gov/tox1/QC/Tox21_112695 <l< td=""></l<>

Chemical "Universe" problem



Biodegradation products Metabolome Protein & DNA fragments Virtual screening libraries Combinatorial chemistry Polymer fragments Polyorganic acids Adducts, surfactants

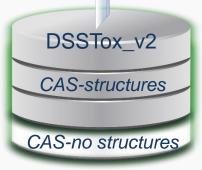
The other 98%

Scientific literature Toxicology studies Environment/Industry Commercially available

Exposure?

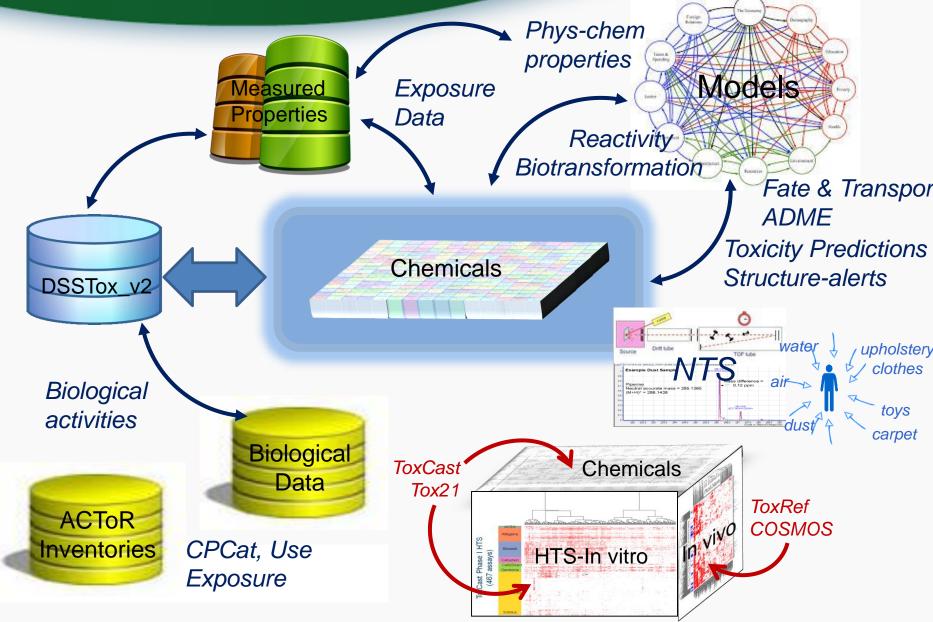
Toxic?

- Where should DSSTox expand chemically?
- What part of the universe should we store in databases?
- How can the valuable ToxCast physical library be shared for greatest gain?
- What cheminformatics "plumbing" would be most useful to this community?



Building EPA's Chemical Informatics Infrastructure & Linkages





Acknowledgements:

- Chris Grulke Lockheed Martin Contractor to EPA DSSTox lead developer/ programmer
- Indira Thillainadarajah EPA SEEP
 DSSTox lead curator
- Kamel Mansouri ORISE Post Doctoral Fellow KNIME workflow
- Jayaram Kancherla ORISE Pre-Doctoral Fellow Chemical dashboard, web-tools development
- Richard Judson NCCT
 ACToR Project Lead
- Antony Williams NCCT
 Cheminformatics Lead
- Mark Strynar, Jon Sobus, Julia Rager, John Wambaugh EPA

This work was reviewed by EPA and approved for publication but does not necessarily reflect official Agency policy.