# 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*
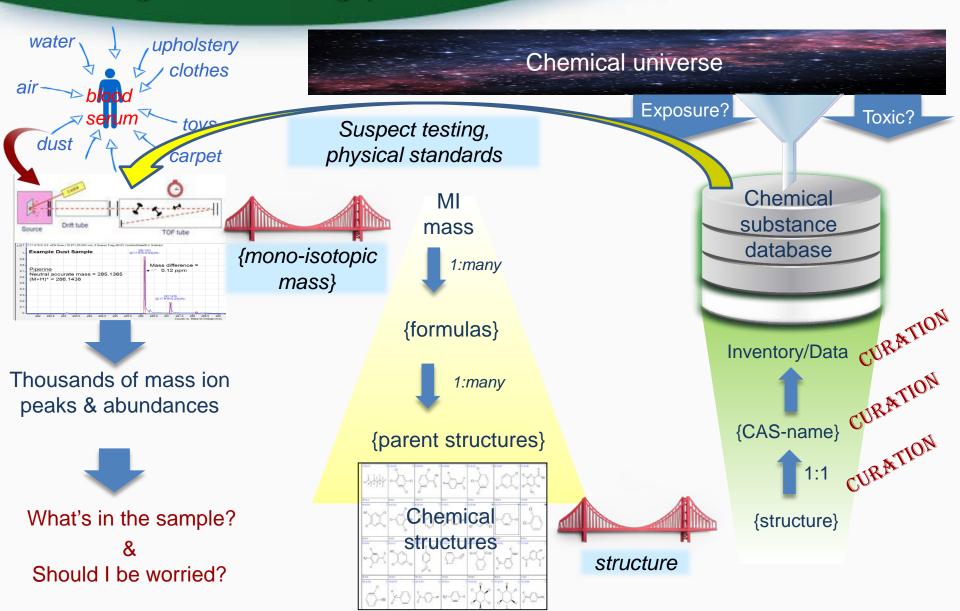
*August 18-19, 2015*
*Research Triangle Park, NC*

0

# Outline

❑ Cheminformatics view of problem

❑ DSSTox chemical database

❑ ToxCast chemical library

❑ Tox21 analytical QC

❑ Challenges

# Cheminformatics view of non-targeted testing problem



EPA
United States
Environmental Protection
Agency

Chemical universe

water
upholstery
clothes
air
blood serum
dust
toys
carpet

Exposure?
Toxic?

Suspect testing, physical standards

Source | Drift tube | TOF tube

Example Dust Sample

Piperine
Neutral accurate mass = 285.1365
(M+H)⁺ = 286.1438

Mass difference = 0.12 ppm

*{mono-isotopic mass}*

MI mass

*1:many*

{formulas}

*1:many*

{parent structures}

Chemical structures

*structure*

Chemical substance database

Inventory/Data

CURATION

{CAS-name}

CURATION

1:1

CURATION

{structure}

Thousands of mass ion peaks & abundances

What's in the sample?
&
Should I be worried?
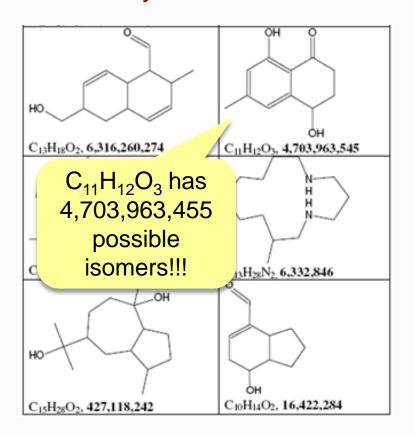
# How big is the problem?

**How Many Structures Can You Generate From A Molecular Formula?**

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

Copyright©2008 Antony Williams

*Highest frequency formula ChemSpider?*

*Theoretically…*



$C_{13}H_{18}O_2$, 6,316,260,274
$C_{11}H_{12}O_3$, 4,703,963,545

$C_{11}H_{12}O_3$ has 4,703,963,455 possible isomers!!!

$C_{15}H_{28}O_2$, 427,118,242
$C_{10}H_{14}O_2$, 16,422,284

| ID | Structure | Empirical Formula | Molecular Weight | Monoisotopic Mass, Da |
|---|---|---|---|---|
| 33808 | | $C_{18}H_{20}N_2O_3$ | 312.363 | 312.147393 |
| 81525 | | $C_{18}H_{20}N_2O_3$ | | |
| 103086 | | $C_{18}H_{20}N_2O_3$ | 312.363 | 312.147393 |

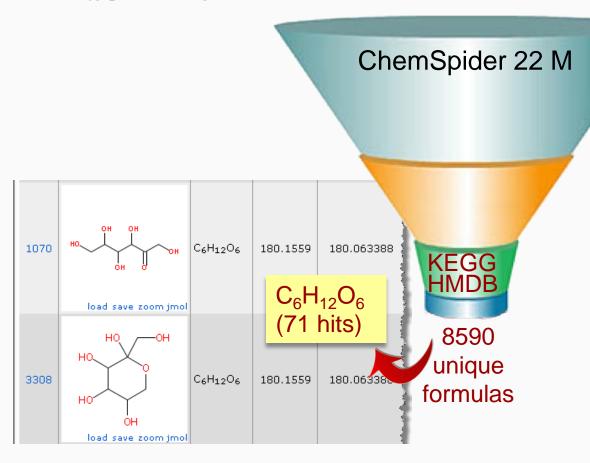$C_{18}H_{20}N_2O_3$ occurs 5110 times

# How big is the problem?

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

**How Many Structures Can You Generate From A Molecular Formula?**
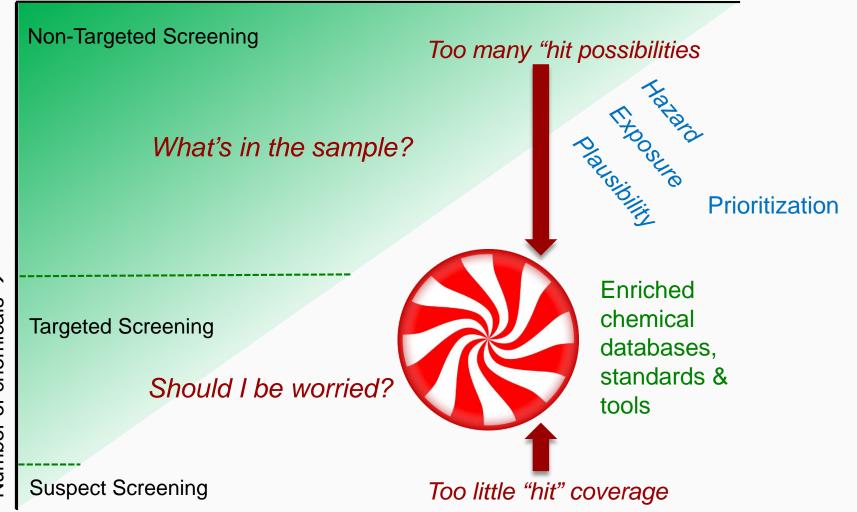Posted by: Antony Williams in ChemSpider Chemistry    *(22 million ChemSpider IDs in 2008)*

Copyright©2008 Antony Williams

ChemSpider 22 M

| 1070 | HO, OH OH / HO OH O OH<br>load save zoom jmol | $C_6H_{12}O_6$ | 180.1559 | 180.063388 |
|---|---|---|---|---|
| 3308 | HO OH / HO O / HO OH<br>load save zoom jmol | $C_6H_{12}O_6$ | 180.1559 | 180.063388 |

$C_6H_{12}O_6$
(71 hits)

KEGG
HMDB

8590
unique
formulas

Problem too big to solve
with formula matching &
structure enumeration

# Cheminformatics view of non-targeted testing problem

# DSSTox_v1 (thru 3/2014)

- 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



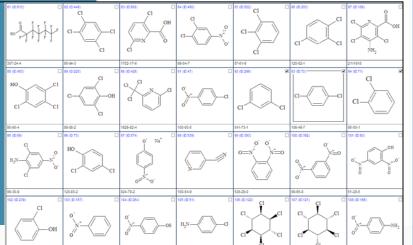**National Center for Computational Toxicology (NCCT)**

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 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

- 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
- TOX21S_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_DefinedOrganic |
| STRUCTURE_ChemicalName_IUPAC |
| STRUCTURE_SMILES |
| STRUCTURE_Parent_SMILES |
| STRUCTURE_InChIS |
| STRUCTURE_InChIKey |
| Substance_modify_yyyymmdd |

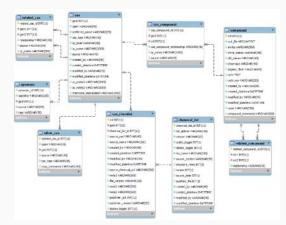*Approx. 25K CAS-substances, 16K structures*

# DSSTox Update

## DSSTox_v1



> 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
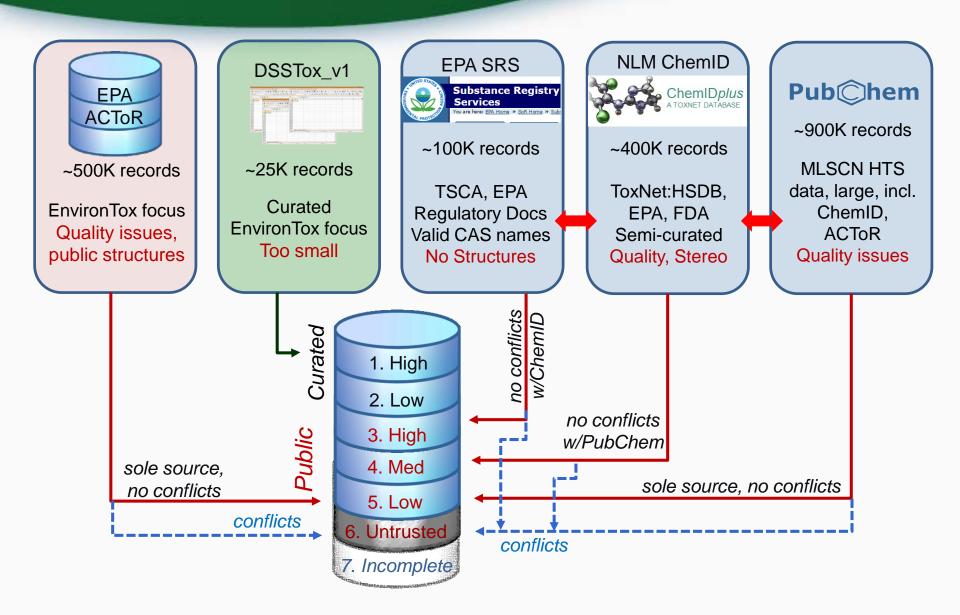
## DSSTox_v2

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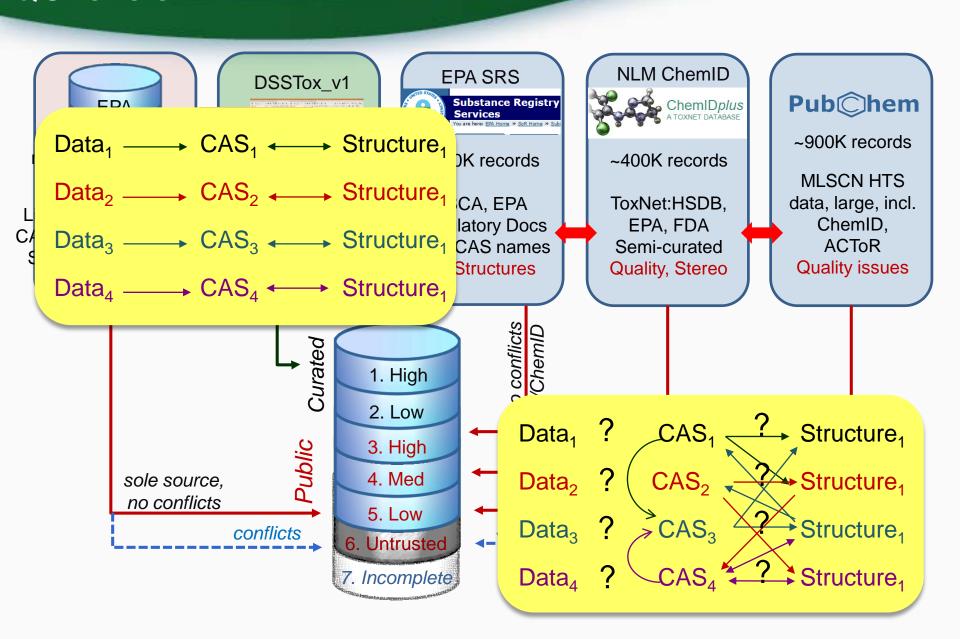> Web-services & Dashboard access

# 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)
   - ✓ *systematic name→ 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"

$CAS_2$ ?

$CAS_1$ ?

Data$_1$ → • Deleted CAS
Data$_2$ → • Invalid CAS
Data$_3$ → • Salt forms
Data$_4$ → • Complex forms
Data$_5$ → • Hydrate forms
Data$_6$ → • Approx mappings to mixtures
Data$_7$ → • Approx mappings to ill-defined substances
Data$_8$ →
Data$_9$ → • Stereoisomers
• Unresolved tautomers

$CAS_5$ ?

NOCAS?   $CAS_4$ ?

$CAS_3$ ?

*DSSTox_v2 Database
& Cheminformatics Layer*

*many:1*

Monoisotopic Mass

↓

Formula

↓

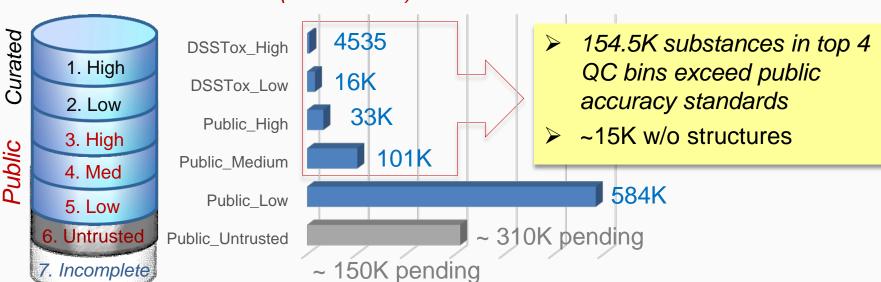Parent structure
(no stereo, desalted)

↓

Valid CAS-substance?

➢ Resolve CAS-structure mappings for accurate chemical-data mapping, i.e. what was tested?

➢ Collapse sphere to bring all related data to NTS parent structure-formula level

# DSSTox_v2 Totals

## QC Level Totals (12Jun2015)

Curated

Public

1. High
2. Low
3. High
4. Med
5. Low
6. Untrusted
7. *Incomplete*

| | |
|---|---|
| DSSTox_High | 4535 |
| DSSTox_Low | 16K |
| Public_High | 33K |
| Public_Medium | 101K |
| Public_Low | 584K |
| Public_Untrusted | ~ 310K pending |

~ 150K pending

> *154.5K substances in top 4 QC bins exceed public accuracy standards*
> ~15K w/o structures

### QC Levels

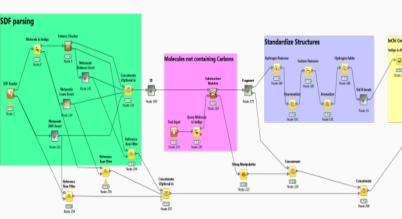| | |
|---|---|
| DSSTox_High: | Hand curated and validated |
| DSSTox_Low: | Hand curated and confirmed using multiple public sources |
| Public_High: | Extracted from EPA SRS and confirmed to have no conflicts in ChemID and PubChem |
| Public_Medium: | Extracted from ChemID and confirmed to have no 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



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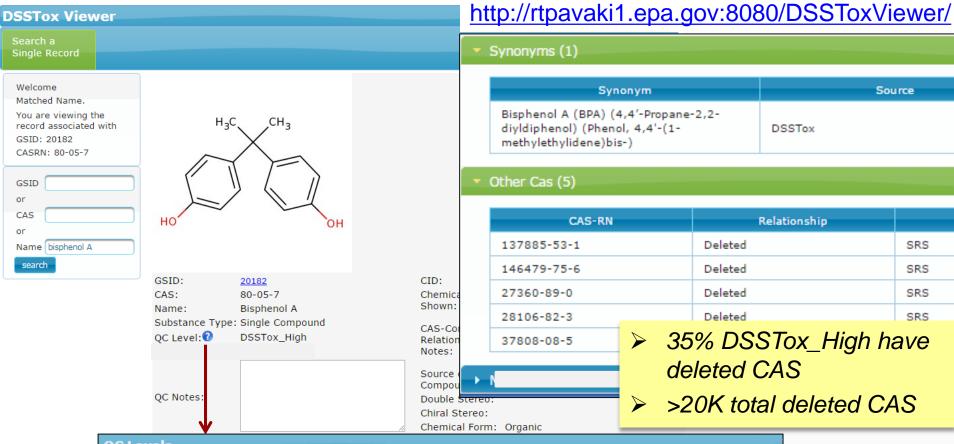✓ Neutralize (when possible)

*Publicly available cheminformatics toolkits in KNIME:*

**Indigo**

RDKit
Open-Source Cheminformatics
and Machine Learning

Pharmaceutical Data
Exploration Laboratory

*Slide courtesy of K. Mansouri*

14

# DSSTox Viewer (EPA Intranet)

**⟆EPA** United States Environmental Protection Agency

## DSSTox Viewer

**Search a Single Record**

Welcome Matched Name.
You are viewing the record associated with
GSID: 20182
CASRN: 80-05-7

GSID [ ]
or
CAS [ ]
or
Name [ bisphenol A ]

[ search ]

H₃C — CH₃ structure (Bisphenol A): HO—C₆H₄—C(CH₃)₂—C₆H₄—OH

GSID: **20182**
CAS: 80-05-7
Name: Bisphenol A
Substance Type: Single Compound
QC Level: ❓ DSSTox_High

QC Notes:

CID:
Chemical Shown:
CAS-Compound Relation Notes:
Source Compound
Double Stereo:
Chiral Stereo:
Chemical Form: Organic

### ▾ Synonyms (1)

| Synonym | Source |
|---|---|
| Bisphenol A (BPA) (4,4'-Propane-2,2-diyldiphenol) (Phenol, 4,4'-(1-methylethylidene)bis-) | DSSTox |

### ▾ Other Cas (5)

| CAS-RN | Relationship | |
|---|---|---|
| 137885-53-1 | Deleted | SRS |
| 146479-75-6 | Deleted | SRS |
| 27360-89-0 | Deleted | SRS |
| 28106-82-3 | Deleted | SRS |
| 37808-08-5 | | |

> ➤ *35% DSSTox_High have deleted CAS*
>
> ➤ *>20K total deleted CAS*

### QC Levels

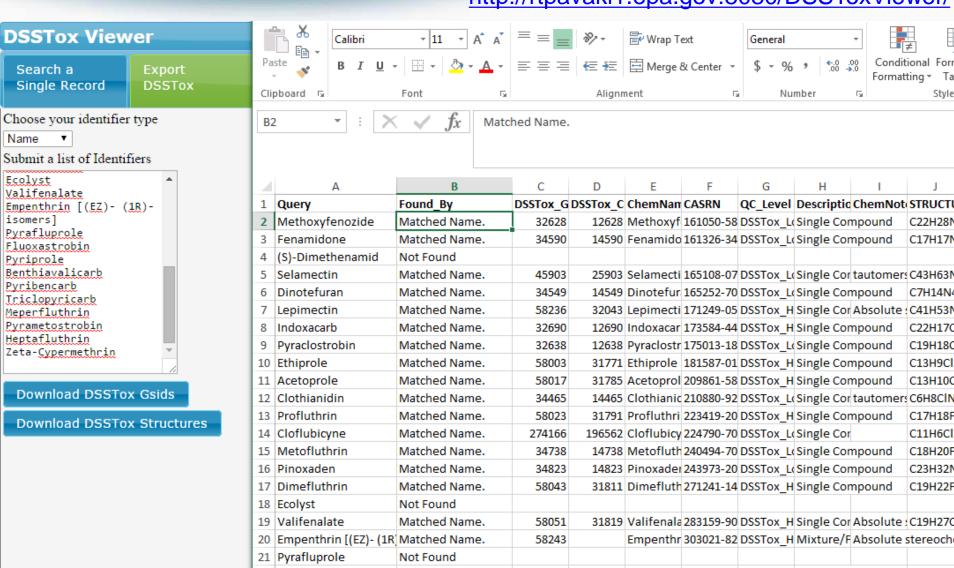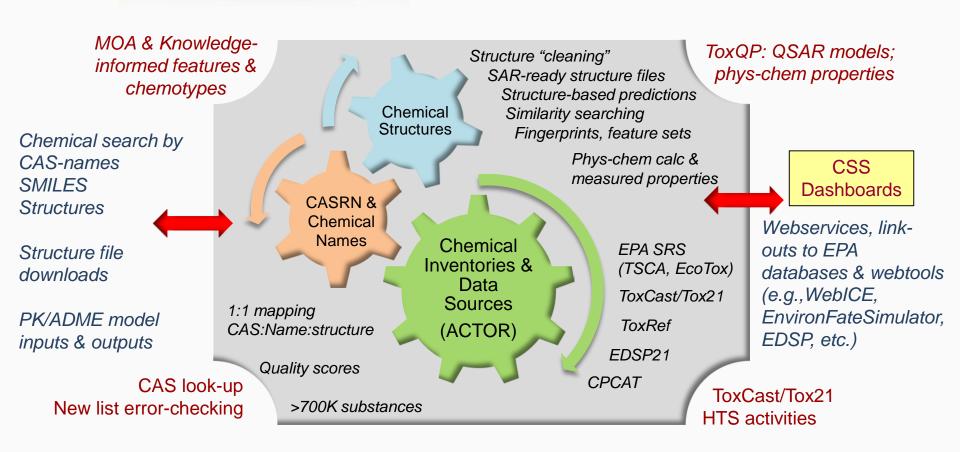| | |
|---|---|
| DSSTox_High: | Hand curated and validated |
| DSSTox_Low: | Hand curated and confirmed using multiple public sources |
| Public_High: | Extracted from EPA SRS and confirmed to have no conflicts in ChemID and PubChem |
| Public_Medium: | Extracted from ChemID and confirmed to have no conflicts in PubChem |
| Public_Low: | Extracted from ACToR or PubChem |
| Public_Untrusted: | Postulated, but found to have conflicts in public sources |

# DSSTox Batch Tool (EPA Intranet)

http://rtpavaki1.epa.gov:8080/DSSToxViewer/
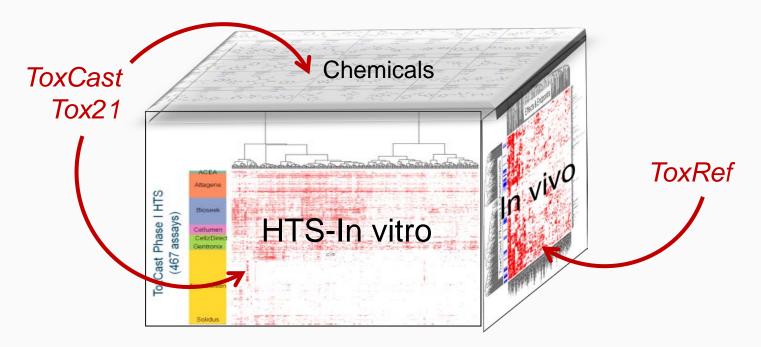
# CSS Chemical Explorer Dashboard
## (powered by DSSTox & ACToR)

*MOA & Knowledge-informed features & chemotypes*

*Structure "cleaning"*
*SAR-ready structure files*
*Structure-based predictions*
*Similarity searching*
*Fingerprints, feature sets*

*ToxQP: QSAR models; phys-chem properties*

Chemical Structures

*Phys-chem calc & measured properties*

*Chemical search by CAS-names SMILES Structures*

*Structure file downloads*

*PK/ADME model inputs & outputs*

CASRN & Chemical Names

Chemical Inventories & Data Sources (ACTOR)

CSS Dashboards

*Webservices, link-outs to EPA databases & webtools (e.g.,WebICE, EnvironFateSimulator, EDSP, etc.)*

*EPA SRS (TSCA, EcoTox)*

*ToxCast/Tox21*

*ToxRef*

*EDSP21*

*CPCAT*

*1:1 mapping CAS:Name:structure*

*Quality scores*

*>700K substances*

CAS look-up
New list error-checking

ToxCast/Tox21
HTS activities
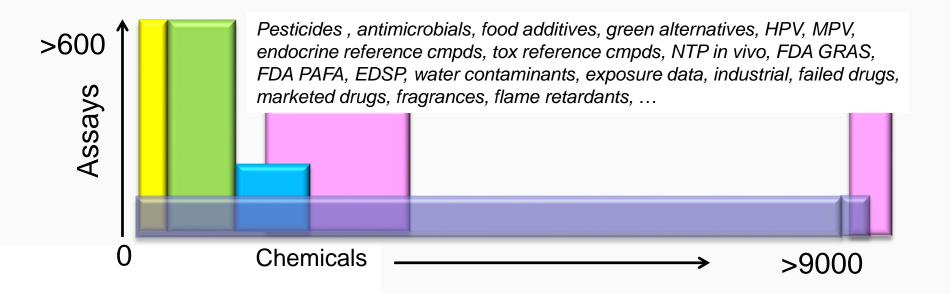
# EPA's ToxCast/Tox21 Projects

- 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
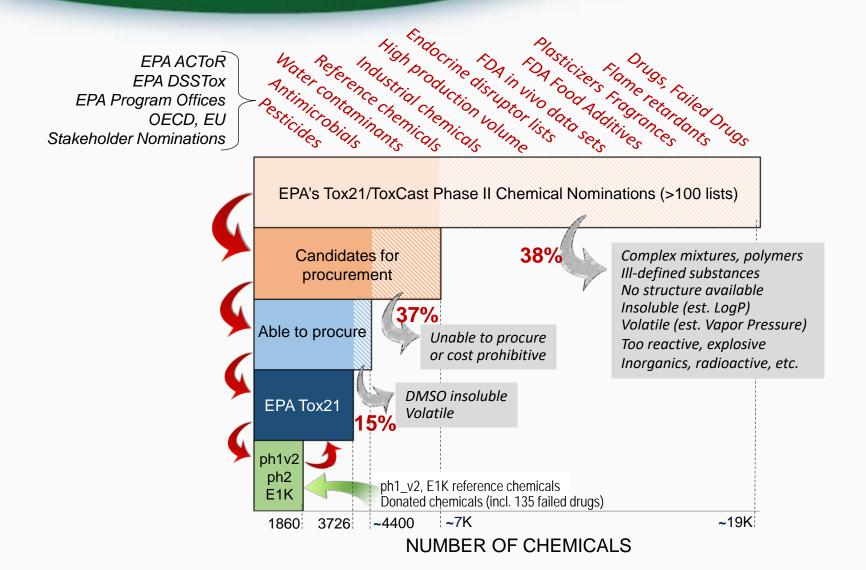
# ToxCast & Tox21 Inventories

| 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 |

*Pesticides , antimicrobials, food additives, green alternatives, HPV, MPV, endocrine reference cmpds, tox reference cmpds, NTP in vivo, FDA GRAS, FDA PAFA, EDSP, water contaminants, exposure data, industrial, failed drugs, marketed drugs, fragrances, flame retardants, …*

>600

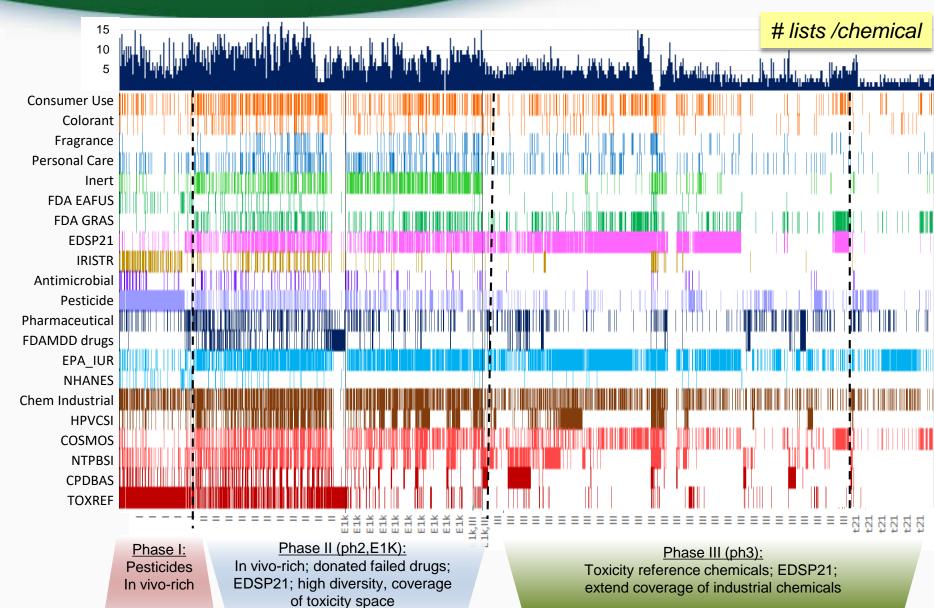Assays
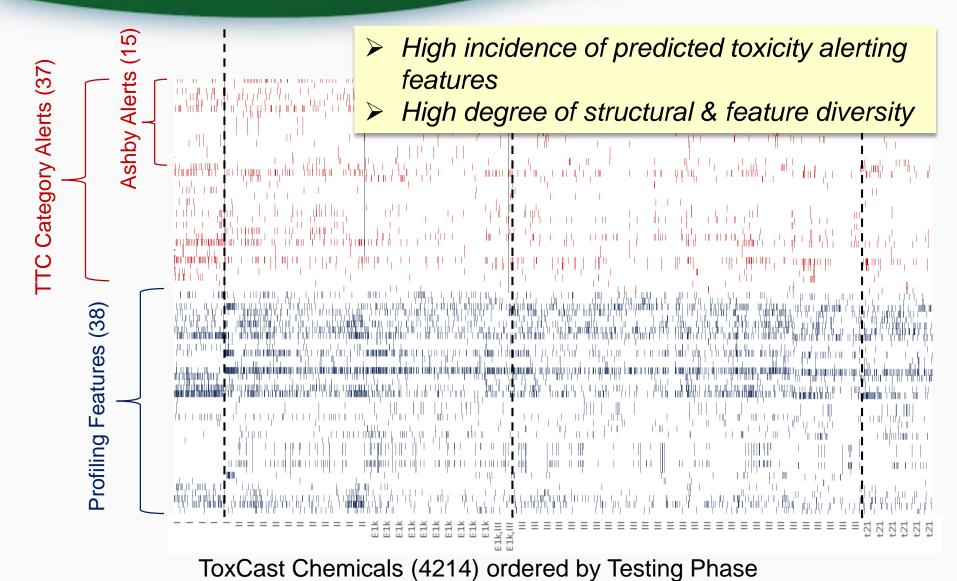
0   Chemicals   →   >9000

# Construction of EPA's Tox21/ToxCast Inventories

ToxCast Chemical Coverage: Use, Exposure, Toxicity

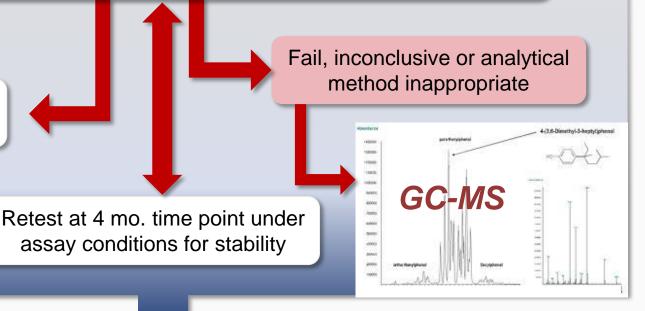# ToxCast: Toxicity Structure-Alerts & Generic Feature Coverage



> ➤ *High incidence of predicted toxicity alerting features*
> ➤ *High degree of structural & feature diversity*

ToxCast Chemicals (4214) ordered by Testing Phase

# Tox21 Analytical QC

**LC-MS**

A copy of each parent Tox21 384 well plate is subjected to analytical QC for assessing purity, identity, concentration, stability

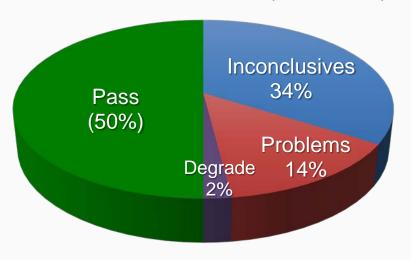Confirm parent ion peak and >90% purity

Fail, inconclusive or analytical method inappropriate

Retest at 4 mo. time point under assay conditions for stability


**GC-MS**

*Work performed under NIH Contract with OpAns, Durham, NC*

Publish QC summary results in association with assay data

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

Tox21_QC_Sum-GSID (8593 total)



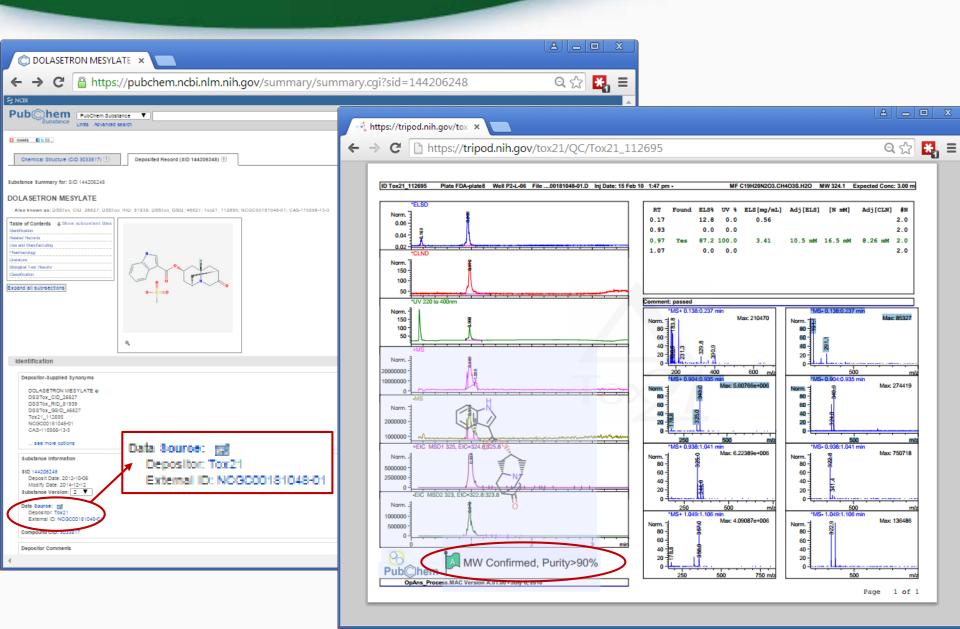Pass (50%)
Inconclusives 34%
Problems 14%
Degrade 2%

- 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?*

*Results as of 8/2015*

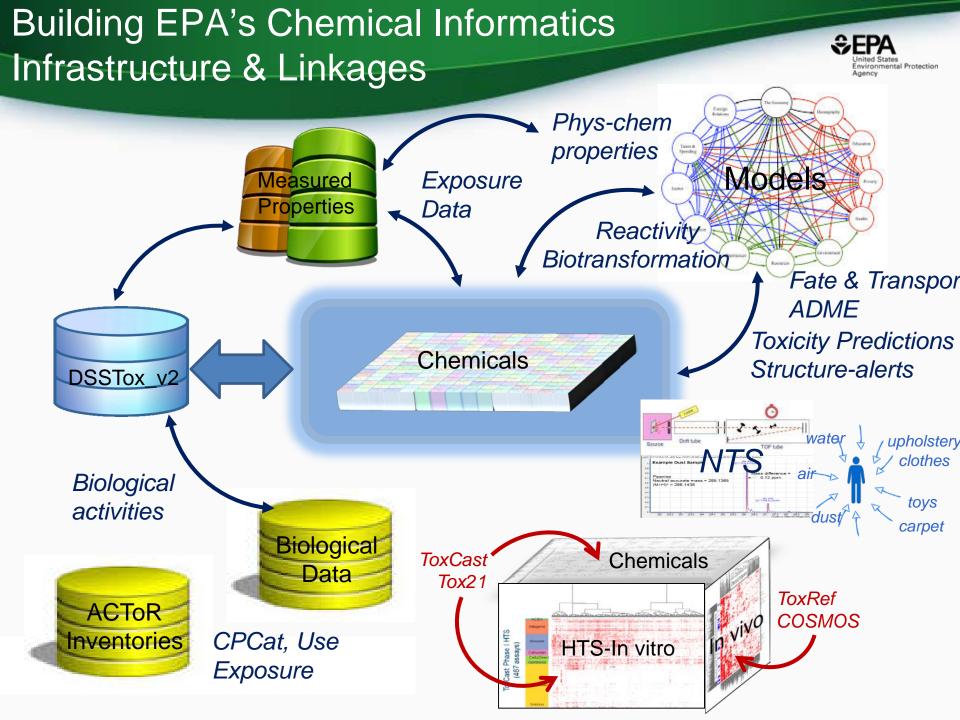# Tox21 Analytical Chemical QC:
## *Publicly available in PubChem*

# 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?

DSSTox_v2

*CAS-structures*

*CAS-no structures*

- ➢ 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