EPA
United States
Environmental Protection
Agency

# The EPA CompTox Dashboard as a Data Integration Hub for Environmental Chemistry Data

*Antony Williams[1], Andrew McEachran[2], Imran Shah[1], Richard Judson[1], John Wambaugh[1], Nancy Baker[3], George Helman[2], Chris Grulke[1], Kamel Mansouri[4], Grace Patlewicz[1], Ann Richard[1], Jeremy Dunne[1] and Jeff Edwards[1]*

1) *National Center for Computational Toxicology, U.S. Environmental Protection Agency, RTP, NC*
2) *Oak Ridge Institute of Science and Education (ORISE) Research Participant, RTP, NC*
3) *Leidos, RTP, NC*
4) *Integrated Laboratory Systems, Inc., RTP, NC*

*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*

*August 2018*
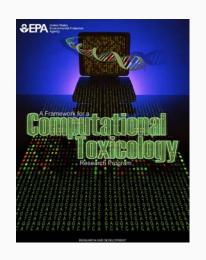*ACS Fall Meeting, Boston*

# Abstract

- The US EPA's CompTox Chemistry Dashboard provides access to various types of data associated with ~760,000 chemical substances. These data include experimental and predicted property data, high-throughput screening assay data and hazard and environmental exposure data. With millions of individual data points and annotations associated with hundreds of thousands of chemicals, data quality is a priority. With tens of thousands of individual users per month browsing the data on the dashboard, the ability of users to provide feedback has allowed us to identify, confirm and address issues in the data. This has required the implementation of novel approaches for data feedback via the user interface that include general feedback on the dashboard and down to individual data points contained in a table. We are presently investigating ways to garner feedback on our ToxCast bioassay data to facilitate the curation of tens of thousands of data points. This presentation will provide an overview of our existing capabilities in the CompTox Chemistry Dashboard for gathering crowdsourced data from the user base and its impact on assisting in the curation of data.

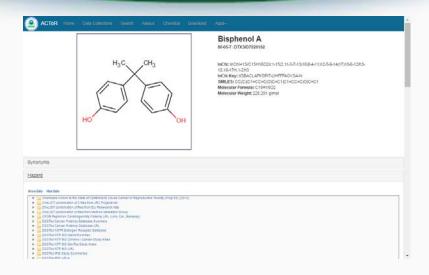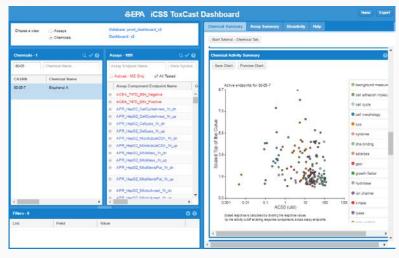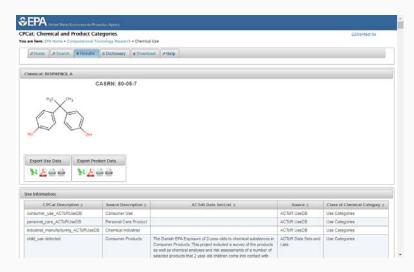# National Center for Computational Toxicology

- National Center for Computational Toxicology established in 2005 to integrate:
  - High-throughput and high-content technologies
  - Modern molecular biology
  - Data mining and statistical modeling
  - Computational biology and chemistry
- Researching computational approaches to quickly evaluate the safety of chemicals for potential risk.
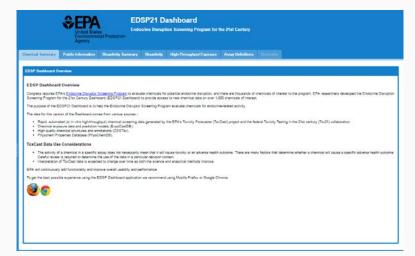- Outputs: a lot of data, models, algorithms and software applications

# Early Dashboard Applications

# The CompTox ~~Chemistry~~ Dashboard

- A publicly accessible website delivering access:
    - New entry portal for all NCCT dashboards
    - ~**762,000** chemicals with related property data
    - **Searchable by chemical, product use, gene and assay (ToxCast)**
    - Experimental and predicted physicochemical property data
    - **"Bioactivity data" for the ToxCast/Tox21 project**
    - **Generalized Read-Across (GenRA) module**
    - Links to other agency websites and public data resources
    - "Literature" searches for chemicals using public resources
    - "Batch searching" for thousands of chemicals
    - DOWNLOADABLE Open Data for reuse and repurposing

# CompTox Portal
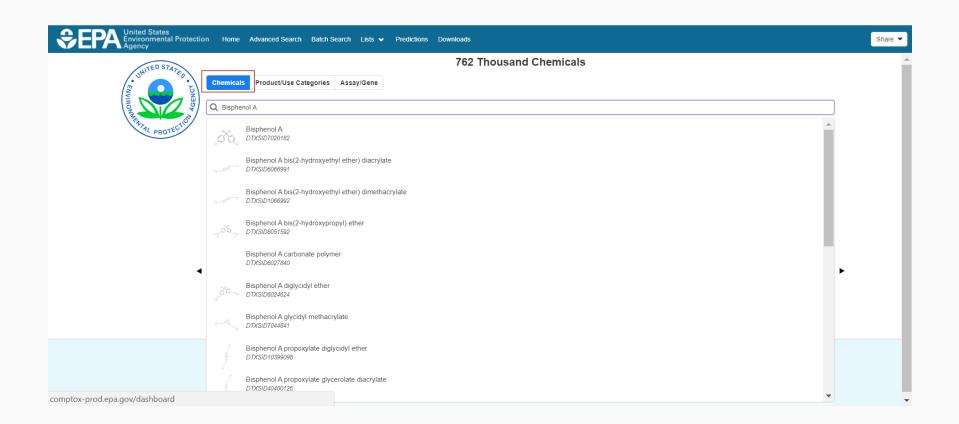
# CompTox Dashboard
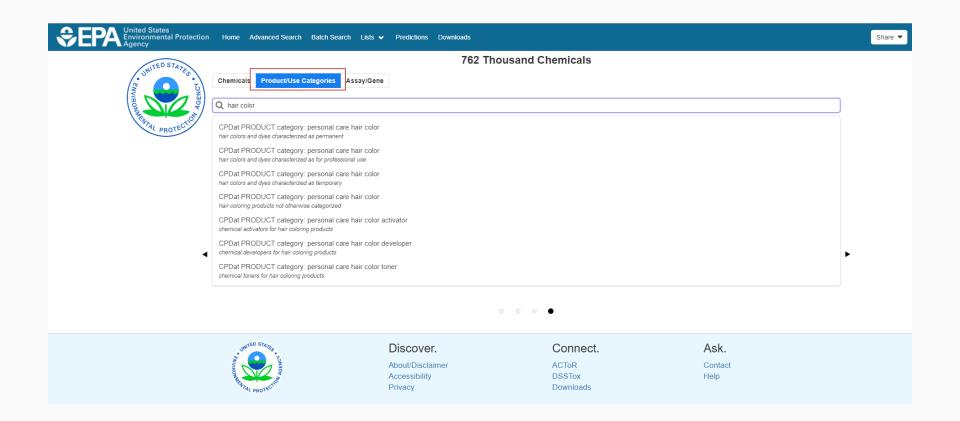https://comptox.epa.gov/dashboard
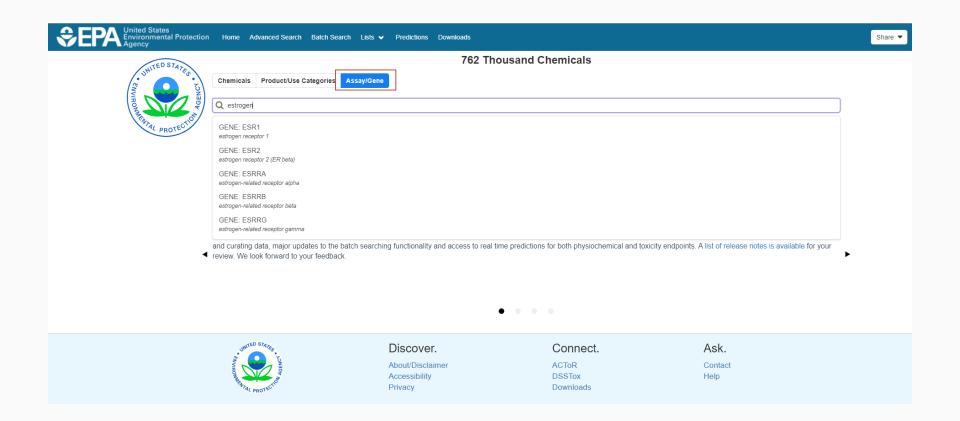
# CompTox Dashboard Chemicals

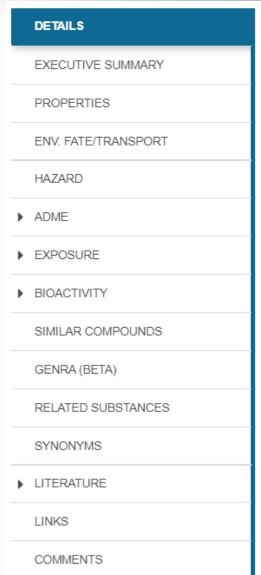# CompTox Dashboard
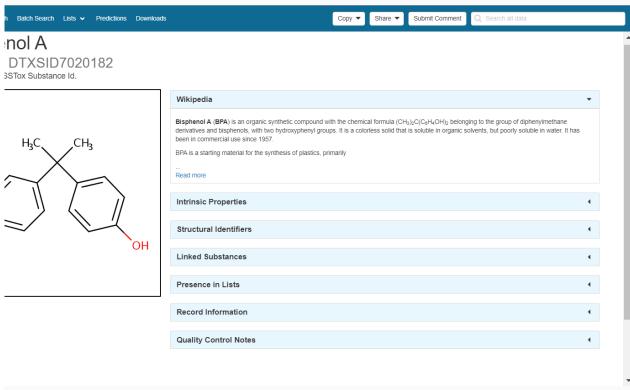# Products and Use Categories

# CompTox Dashboard
## Assays and Genes

# Detailed Chemical Pages

# Physicochemical properties

# OPERA Predicted Properties

**Journal of Cheminformatics**

**RESEARCH ARTICLE**  **Open Access**

CrossMark

# OPERA models for predicting physicochemical properties and environmental fate endpoints

Kamel Mansouri[1,2,3*], Chris M. Grulke[1], Richard S. Judson[1] and Antony J. Williams[1]

# Detailed OPERA Prediction Reports

# Access to Chemical Hazard Data

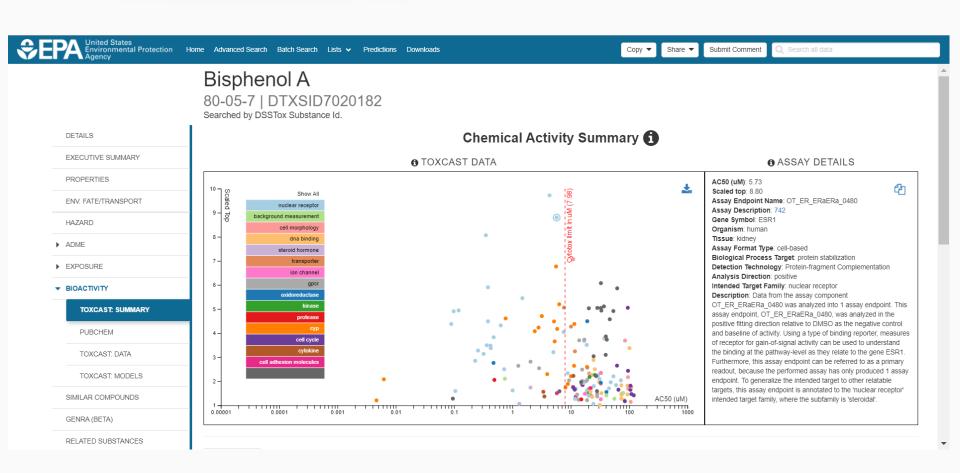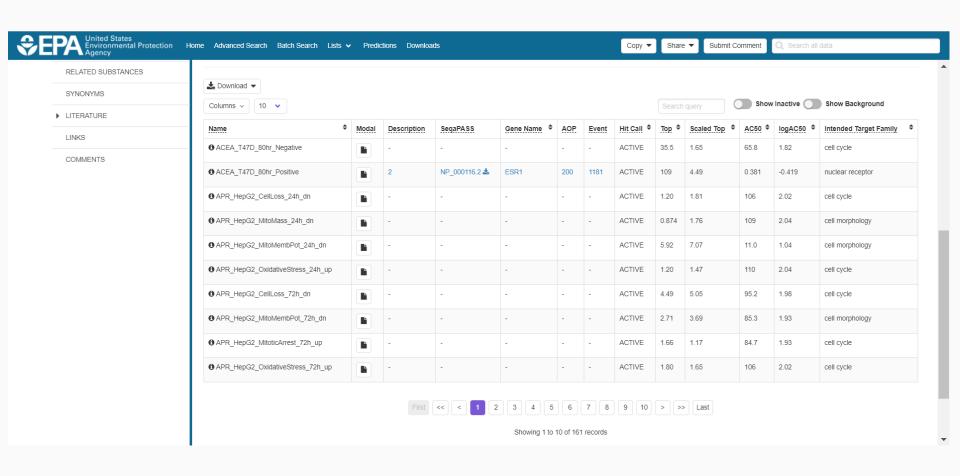# *In Vitro* Bioassay Screening
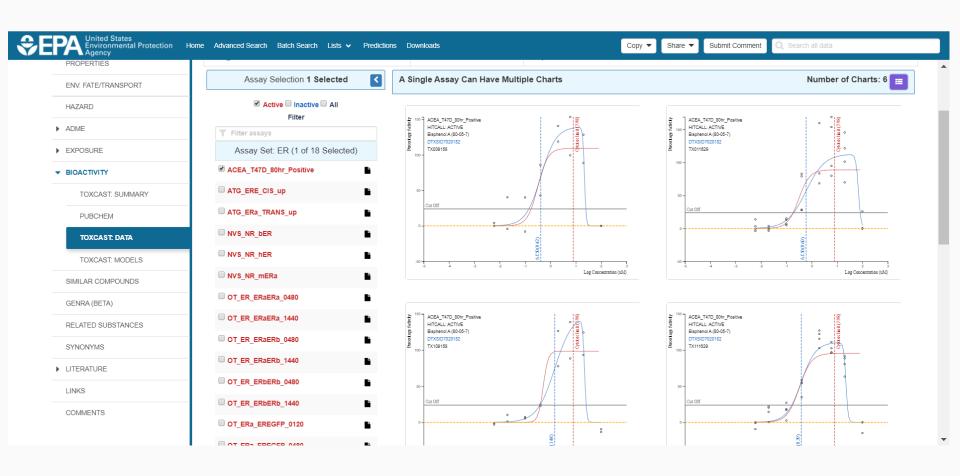## ToxCast and Tox21

# *In Vitro* Bioassay Screening
## ToxCast and Tox21

# *In Vitro* Bioassay Screening
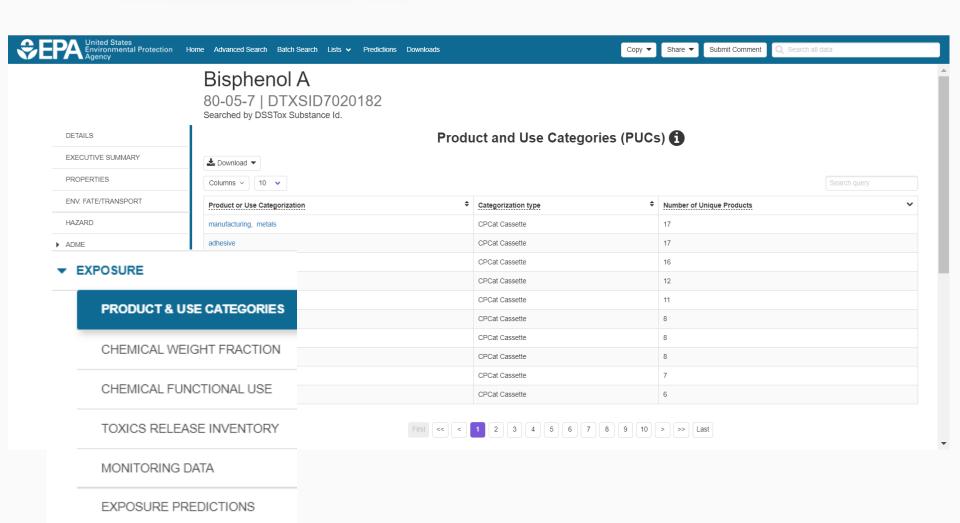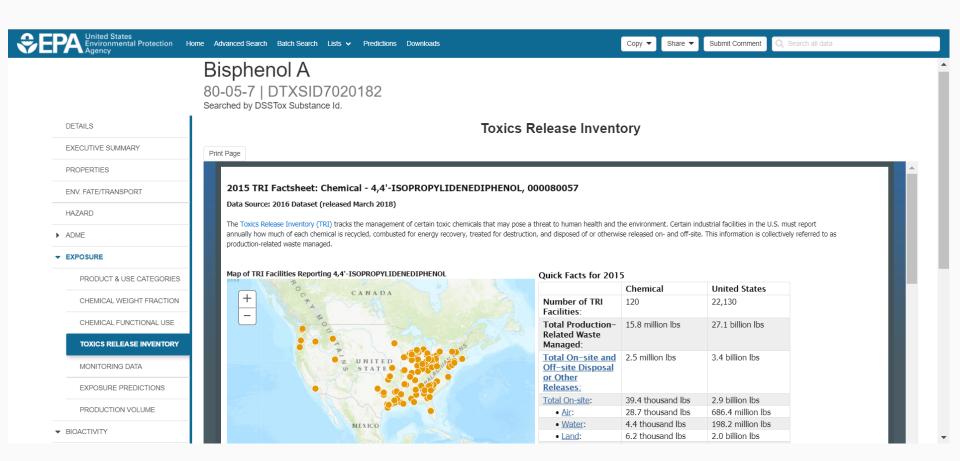## ToxCast and Tox21

# Sources of Exposure to Chemicals

# Sources of Exposure to Chemicals

# Generalized Read-Across (GenRA)

# Identifiers to Support Searches

# Literature Searches and Links

# Abstract Sifter – PubMed Integration

# Abstract Sifter for Excel



SOFTWARE TOOL ARTICLE

## Abstract Sifter: a comprehensive front-end system to PubMed [version 1; referees: 2 approved]

Nancy Baker [1], Thomas Knudsen[2], Antony Williams [2]

+ Author details

This article is included in the Chemical Information Science gateway.

## Abstract

The Abstract Sifter is a Microsoft Excel based application that enhances existing search capabilities of PubMed. The Abstract Sifter assists researchers to search effectively, triage results, and keep track of articles of interest. The tool implements an innovative "sifter" functionality for relevance ranking, giving the researcher a way to find articles of interest quickly. The tool also gives

# External Links to ~80 websites

# Integrated Linkouts

A Substance Registry ...    ACToR    Toxline

eChemPortal provides free public access to information on properties of chemicals. Direct links to collections of chemical hazard and risk information prepared for government chemical review programmes at national, regional and international levels are obtained.

kipedia    eChemPortal    G Google Scholar

DS Lookup    Gene-Tox    G Google Patents

MSDS Lookup    Gene-Tox    G

The International Chemical Safety Cards (ICSC) summarize essential health and safety information on chemicals for their use at the

NIOSH Chemical Safety ...    LactMed    IR

Substance Registry ...    ACToR    Toxline

Comparative Toxicogenomics Database is a robust, publicly available database that aims to advance understanding about how environmental exposures affect human health.

Bank    CTD    G Google Books

pedia    eChemPortal    G Google Scholar

# Integrated Linkouts
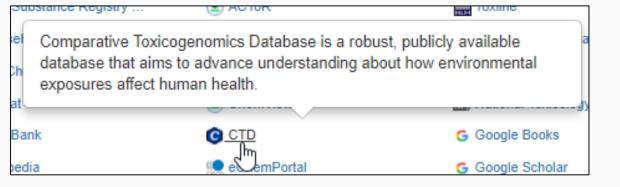## Comparative Toxicogenomics DB

# Not just chemical "structures"

- ## Chemicals in commerce, of interest to the EPA, are not all easily represented by structures

- ## Different chemical substances supported
  - Chemical structures
  - "UVCB chemicals" - Unknown or Variable Composition, Complex Reaction Products and Biological Materials
  - Metabolites and transformation products
  - Homologous series as Markush Structures
  - Curated classes of chemicals

# UVCB Chemicals

# Markush Structures

# Transformation Products

# Not just chemical "structures"

- Different chemical substances supported
  - Chemical structures
  - "UVCB chemicals" - Unknown or Variable Composition, Complex Reaction Products and Biological Materials
  - Metabolites and transformation products
  - Homologous series as Markush Structures
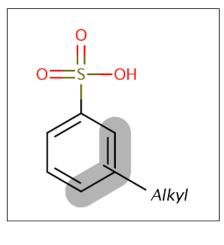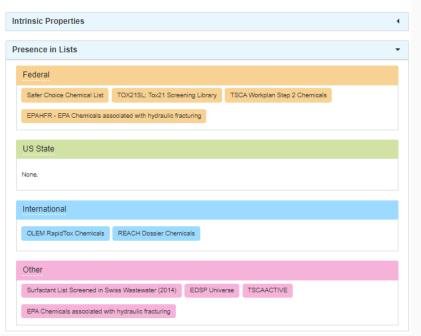  - Curated classes of chemicals

- Lists of chemicals
  - Submitted lists of chemicals – Federal, State, International and other general lists
  - Growing lists to support specific projects – e.g. ToxCast phases, algal toxins, our publication datasets

# List of Chemicals

# Algal Toxins

# Mass and Formula Searches Supporting Mass Spectrometry

# Advanced Searches
# Mass Based Search

# Advanced Searches
# Mass Based Search

# Batch Searching

- Singleton searches are useful but we work with thousands of chemicals!

- Typical questions
  - What are the SMILES strings for a list of 1000 chemicals?
  - Do any of this list of chemicals have XXX type of data?
  - What are the predicted logP values for a list of chemicals?
  - Can I get chemical lists in Excel files? In SDF files?

# Batch Searching



39

# Batch Searching

**Select Output Format:**

| ▦ Excel ⌄ |

**⬇ Download**

**Customize Results**
- ☐ Select All
- ☐ Select All in Lists

**Chemical Identifiers**
- ☑ DTXSID ❶
- ☑ Chemical Name ❶
- ☐ CAS-RN ❶
- ☐ InChIKey ❶
- ☐ IUPAC Name ❶

**Structures**
- ☐ Mol File ❶
- ☐ SMILES ❶
- ☐ InChI String ❶
- ☐ MS-Ready SMILES ❶
- ☐ QSAR-Ready SMILES ❶

**Intrinsic And Predicted Properties**
- ☐ Molecular Formula ❶
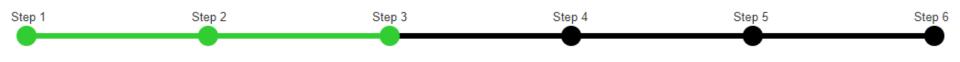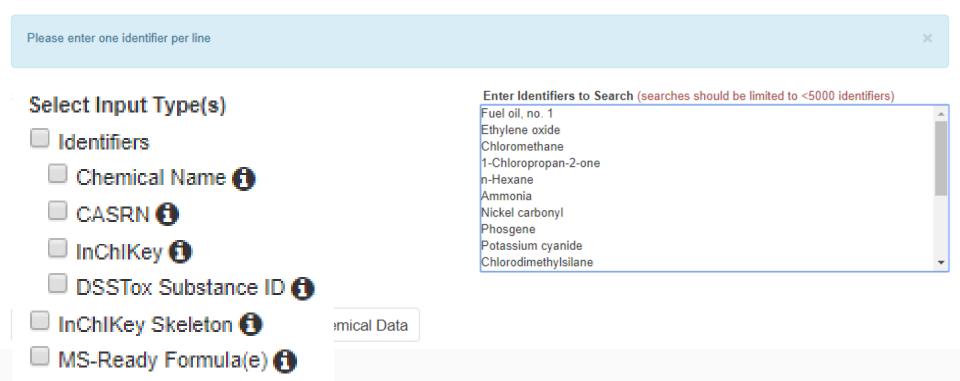- ☐ Average Mass ❶
- ☐ Monoisotopic Mass ❶
- ☐ TEST Model Predictions ❶
- ☐ OPERA Model Predictions ❶

**Presence in Lists:**
- ☐ ICCVAM test method evaluation report: in vitro ocular toxicity test methods
- ☐ 40CFR355
- ☐ A list of all PBDEs (Polybrominated diphenyl ethers)
- ☐ A list of all PCBs (Polychlorinated biphenyls)
- ☐ A list of polycyclic aromatic hydrocarbons
- ☐ Acute exposure guideline levels
- ☐ Algal Toxins
- ☐ Androgen Receptor Chemicals
- ☐ APCRA Chemicals for Prospective Analysis
- ☐ APCRA Chemicals for Retrospective Analysis
- ☐ APCRA Chemicals for Retrospective Analysis_App_List_448_Chemicals
- ☐ ATSDR Minimal Risk Levels (MRLs) for Hazardous Substances
- ☐ ATSDR Toxic Substances Portal Chemical List
- ☐ Bisphenol Compounds
- ☐ California Office of Environmental Health Hazard Assessment
- ☐ Chemicals with interesting names
- ☐ CMAP
- ☐ DNT Screening Library
- ☐ Drinking Water Suspects, KWR Water, Netherlands
- ☐ EDSP Universe
- ☐ EPA Chemicals associated with hydraulic fracturing

# Excel Output
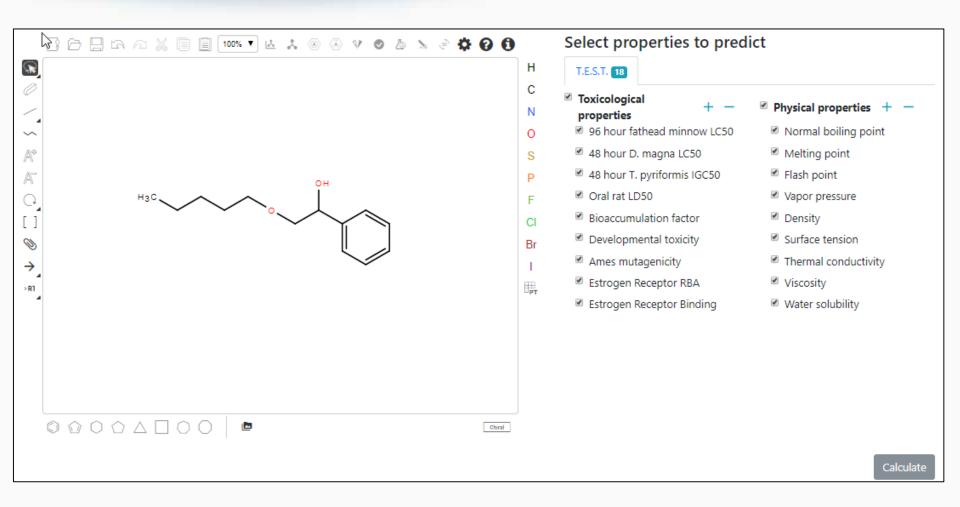
| INPUT | FOUND_BY | DTXCID_IN | DATA_SOL | TOXVAL_D | TOXCAST_ | TOXCAST_ | NUMBER_( | PUBCHEM_ | STO |
|-------|----------|-----------|----------|----------|----------|----------|----------|----------|-----|
| C6H12O3 | MS Ready I | DTXCID701! | 91 | Y | 0.36 | 2/562 | 24 | 83 | Y |
| C6H12O3 | MS Ready I | DTXCID003 | 67 | Y | 0.36 | 1/276 | 376 | 80 | Y |
| C6H12O3 | MS Ready I | DTXCID106 | 65 | Y | 4.42 | 5/113 | 6 | 77 | Y |
| C6H12O3 | MS Ready I | DTXCID105 | 45 | Y | 0.0 | 0/163 | 3 | 94 | - |
| C6H12O3 | MS Ready I | DTXCID901 | 38 | Y | - | - | 14 | 110 | Y |
| C6H12O3 | MS Ready I | DTXCID402 | 34 | Y | 0.0 | 0/113 | - | 53 | Y |
| C6H12O3 | MS Ready I | DTXCID202! | 31 | Y | - | - | - | 36 | Y |
| C6H12O3 | MS Ready I | DTXCID202 | 30 | - | 2.54 | 7/276 | - | 54 | - |
| C6H12O3 | MS Ready I | DTXCID109 | 26 | Y | - | - | - | 46 | - |
| C6H12O3 | MS Ready I | DTXCID202! | 24 | Y | 0.0 | 0/113 | - | 47 | - |
| C6H12O3 | MS Ready I | DTXCID303( | 22 | Y | - | - | - | 89 | - |
| C6H12O3 | MS Ready I | DTXCID302 | 20 | Y | - | - | 2 | 25 | Y |
| C6H12O3 | MS Ready I | DTXCID407 | 19 | Y | - | - | 12 | 62 | - |
| C6H12O3 | MS Ready I | DTXCID704 | 17 | Y | - | - | - | 64 | - |
| C6H12O3 | MS Ready I | DTXCID704 | 16 | Y | - | - | 3 | 49 | - |

# Real-Time Predictions

# Real-Time Predictions

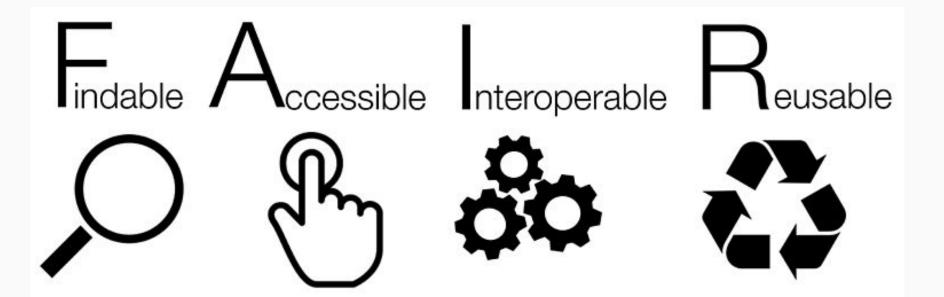| Property | Experimental Value | Prediction | | | | |
|---|---|---|---|---|---|---|
| | | Consensus | Hierarchical clustering | Single model | Group contribution | Nearest neighbor |
| 96 hour fathead minnow LC50 | | 4.477 -Log10(mol/L) 6.954 mg/L | 4.195 -Log10(mol/L) 13.288 mg/L | 3.994 -Log10(mol/L) 21.110 mg/L | 3.478 -Log10(mol/L) 69.224 mg/L | 6.238 -Log10(mol/L) 0.120 mg/L |
| 48 hour D. magna LC50 | | 4.398 -Log10(mol/L) 8.328 mg/L | 3.877 -Log10(mol/L) 27.677 mg/L | 4.039 -Log10(mol/L) 19.026 mg/L | 4.084 -Log10(mol/L) 17.173 mg/L | 5.593 -Log10(mol/L) 0.532 mg/L |
| 48 hour T. pyriformis IGC50 | | 4.063 -Log10(mol/L) 18.039 mg/L | 3.731 -Log10(mol/L) 38.668 mg/L | | 3.386 -Log10(mol/L) 85.610 mg/L | 5.070 -Log10(mol/L) 1.773 mg/L |
| Oral rat LD50 | | 1.758 -Log10(mol/kg) 3640.950 mg/kg | 1.982 -Log10(mol/kg) 2172.756 mg/kg | | | 1.533 -Log10(mol/kg) 6101.245 mg/kg |
| Bioaccumulation factor | | 1.797 Log10 62.700 | 2.202 Log10 159.310 | 1.287 Log10 19.346 | 1.181 Log10 15.157 | 2.520 Log10 330.834 |
| Developmental toxicity | | false | false | false | | true |
| Ames mutagenicity | | false | false | | | false |
| Estrogen Receptor RBA | | -3.075 Log10 $8.418*10^{-4}$ | -3.078 Log10 $8.356*10^{-4}$ | -3.720 Log10 $1.907*10^{-4}$ | | -2.427 Log10 0.004 |
| Estrogen Receptor Binding | | true | true | true | false | true |

# Crowdsourced Curation

# Crowdsourced Curation

- An
vie

Details to be submitted with your comment:

Text selected: 149.999

Found On: August 5th 2018, 10:00:16 pm

Original Query: /dsstoxdb/results?search=BPA#toxicity-values

Browser: Chrome 68

Probably want to round up this value in the Hazard Table.

williams.antony@epa.gov

I'm not a robot

reCAPTCHA
Privacy · Terms

Submit

# Our support for FAIR Data

# Downloadable Data

- Present work in development
  - Real time prediction using OPERA models
  - Structure/substructure/similarity search integration
  - Ongoing expansion of chemicals
  - Release of new ToxCast database (v3_2018)
  - Addition of products data from 10s of thousands of MSDS sheets
  - Analytical Data support
    - Integration of analytical data for ToxCast/Tox21 data
    - Spectral searching against predicted Mass Spectra
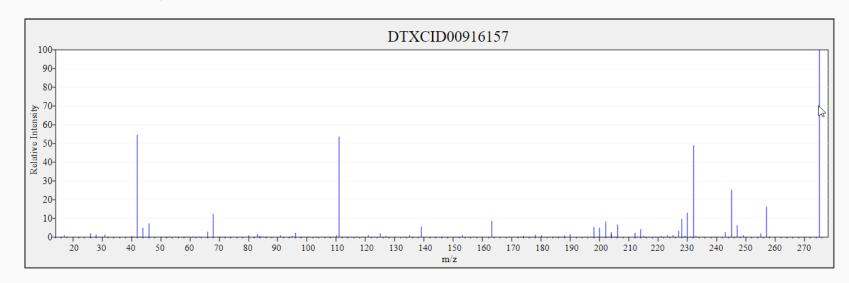
# Prototype Development

# Predicted Mass Spectra

http://cfmid.wishartlab.com/



- MS/MS spectra prediction for ESI+, ESI-, and EI
- Predictions generated and stored for >700,000 structures, to be accessible via Dashboard

# Conclusion

- The EPA CompTox Dashboard provides access to data for ~762,000 chemicals, ToxCast assays and associated product use categories

- High quality data from ongoing curation efforts

- An integration hub for multiple "modules"
  - Experimental and predicted properties
  - Human and Ecological Hazard data
  - Exposure data – products, data in the environment
  - *In vitro* bioassay data – ToxCast/Tox21
  - Literature searching – Google Scholar and PubMed
  - Specialized searches – mass/formula for analytical support
  - Batch searching and Real Time Predictions

- Data and functionality increases with every release

# Acknowledgments

- The NCCT CompTox Chemistry Dashboard Development Team
- NERL scientists (Jon Sobus, Elin Ulrich) – Mass Spectrometry
- Kamel Mansouri – OPERA models
- Todd Martin and Valery Tkachenko – TEST predictions
- Nancy Baker – Abstract Sifter

# Contact

**Antony Williams**

US EPA Office of Research and Development

National Center for Computational Toxicology (NCCT)

Williams.Antony@epa.gov

**ORCID**: https://orcid.org/0000-0002-2668-4821