

# Integration of Markush Structures into EPA's DSSTox Database to Represent and Enumerate UVCB Substances

Christopher Grulke†

Antony Williams

Ann Richard

*NCCT, U.S. EPA*



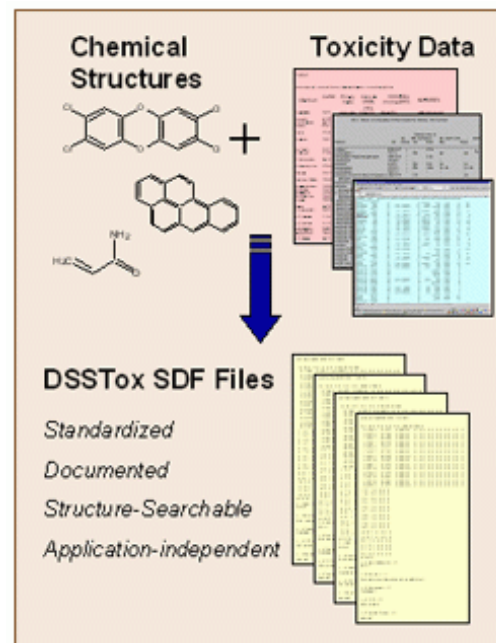
American Chemical Society Meeting, Spring 2018

18 March 2018, New Orleans, LA

# DSSTox Ancient 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 Current History

- 761K substance records (27.5K manually curated)
- Central database for the Comptox Chemical Dashboard
- More Goals:
  - Become a hub for all chemical data relevant to an environmental scientist
  - Provide batch extraction of chemical data for our user community
  - Offer chemical list based views of our data
  - Provide list specific search capabilities
- Check out: <https://comptox.epa.gov/dashboard>

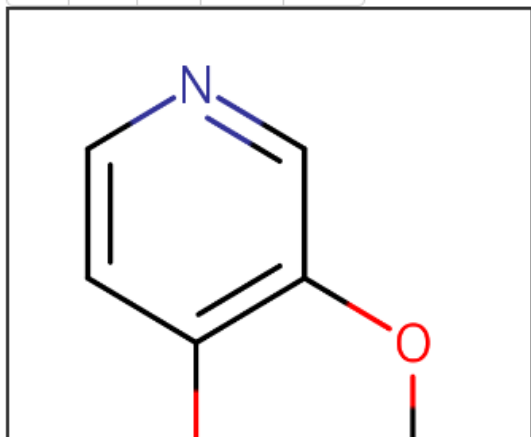
## Chemistry Dashboard

[Submit Comment](#)[Share ▾](#)[Copy ▾](#)[Aa ▾](#)

### 4-Hydroxy-3-methoxypyridine

62885-41-0 | DTXSID80198757

© Searched by CAS-RN: Found 1 result for '62885-41-0'.



#### Intrinsic Properties

**Molecular Formula:** C<sub>6</sub>H<sub>7</sub>NO<sub>2</sub>

[Find All Chemicals](#)



**Average Mass:** 125.127 g/mol



**Monoisotopic Mass:** 125.047678 g/mol



#### Structural Identifiers

#### Related Compounds (Beta)

[About](#)[Contact](#)[Privacy](#)

Powered by **ACToR**



Powered by **DSSTox**

[Accessibility](#)[Help](#)[Downloads](#)

# Comptox Chemistry Dashboard (Cont.)

## Chemistry Dashboard

Submit Comment

Share ▾

Co

Chemical Properties

Env. Fate/Transport

Synonyms

External Links

Toxicity Values (Beta)

Exposure

Bioassays

Similar Molecules (Beta)

Literature

Comments

### Summary

LogP: Octanol-Water

Water Solubility

Density

Melting Point

Boiling Point

Surface Tension

Vapor Pressure

LogKoa: Octanol-Air

Henry's Law

Index of Refraction

Molar Refractivity

nKa Basic Apparent

Download as:

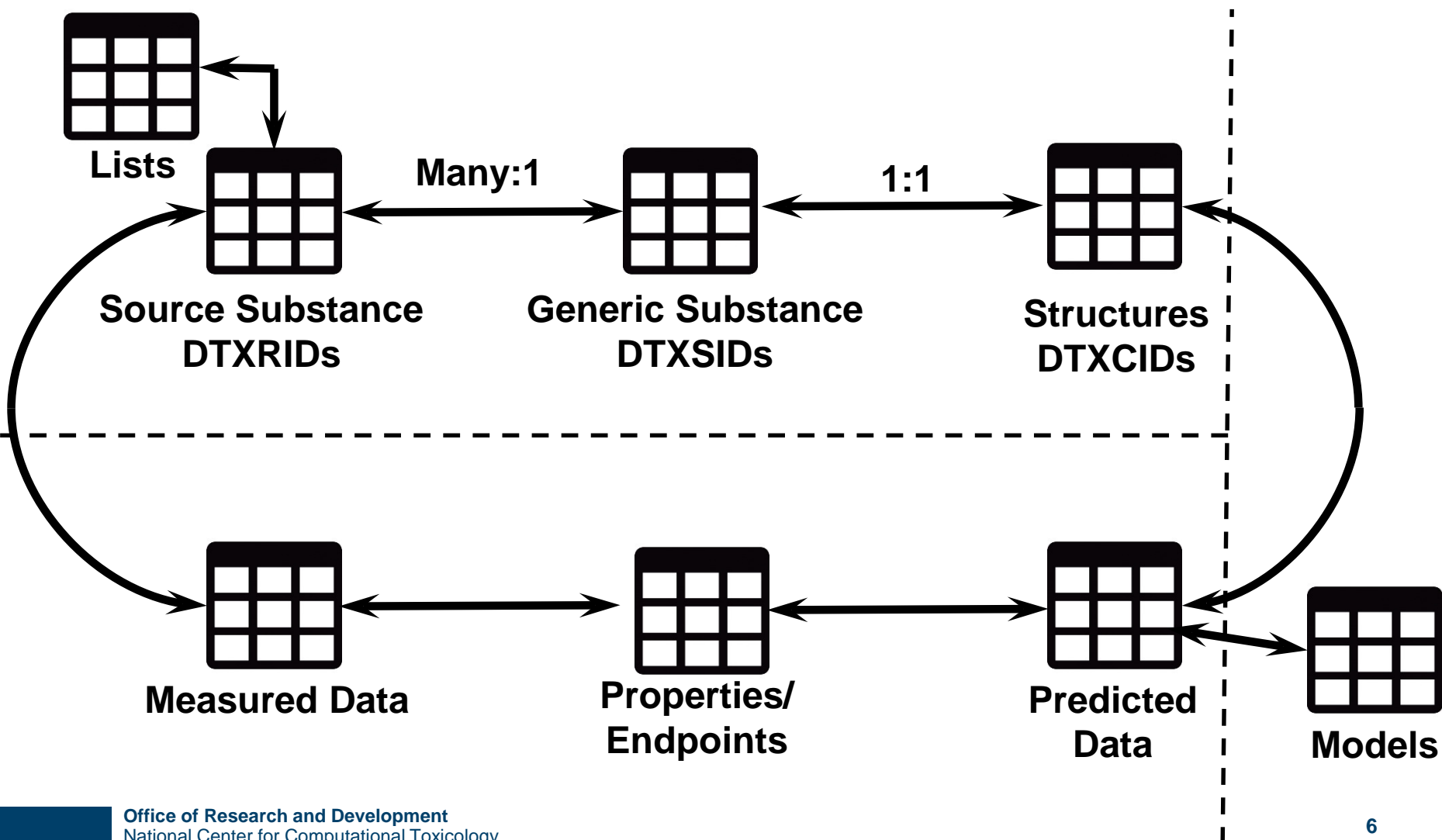
TSV

Excel

SDF

Property	Average		Median		Range		Unit
	Experimental	Predicted	Experimental	Predicted	Experimental	Predicted	
LogP: Octanol-Water	2.61 (1)	2.76 (4)	2.61 to 2.61	2.76	2.61	2.50 to 3.05	-
Water Solubility	1.30e-04 (1)	1.46e-02 (4)	1.30e-04 to 1.30e-04	1.46e-02	1.30e-04	1.50e-04 to 5.71e-02	mol/L
Density	-	1.27 (1)	-	1.27	-	-	g/cm <sup>3</sup>
Melting Point	174 (6)	151 (3)	173 to 177	151	173 to 177	114 to 185	°C
Boiling Point	-	312 (3)	-	312	-	284 to 339	°C
Surface Tension	-	53.8 (1)	-	53.8	-	-	dyn/cm
Vapor Pressure	7.21e-11 (1)	4.47e-06 (3)	7.21e-11 to 7.21e-11	4.47e-06	7.21e-11	2.06e-07 to 1.27e-05	mmHg
LogKoa: Octanol-Air	-	8.38 (1)	-	8.38	-	-	-
Henry's Law	-	4.20e-10 (1)	-	4.20e-10	-	-	atm-m <sup>3</sup> /mol
Index of Refraction	-	1.61 (1)	-	1.61	-	-	-
Molar Refractivity	-	58.5 (1)	-	58.5	-	-	cm <sup>3</sup>
pKa Basic Apparent	-	2.27 (1)	-	2.27	-	-	-
Molar Volume	-	170 (1)	-	170	-	-	cm <sup>3</sup>
Polarizability	-	23.2 (1)	-	23.2	-	-	Å <sup>3</sup>


# General Data Model



# So Where is the Data?

## C10-16 Alcohols

67762-41-8 | DTXSID4028331

 Searched by Approved Name: Found 1 result for 'C10-16 Alcohols'.

### Presence in Lists

Federal

Safer Choice Chemical List

US State

International

Other

TSCAACTIVE

### Record Information

### Quality Control Notes

Related Substances

Synonyms

Links

Bioassays

Exposure

Hazard

Comments

Chemical Properties

Literature

No Chemical Properties Found.

# UVCB Chemicals

[Environmental Topics](#)[Laws & Regulations](#)[About EPA](#)[Search EPA.gov](#)

## TSCA Chemical Substance Inventory

[CONTACT US](#)[SHARE](#)[TSCA Inventory Home](#)[About the Inventory](#)[Access the Inventory](#)[Policy and Guidance](#)

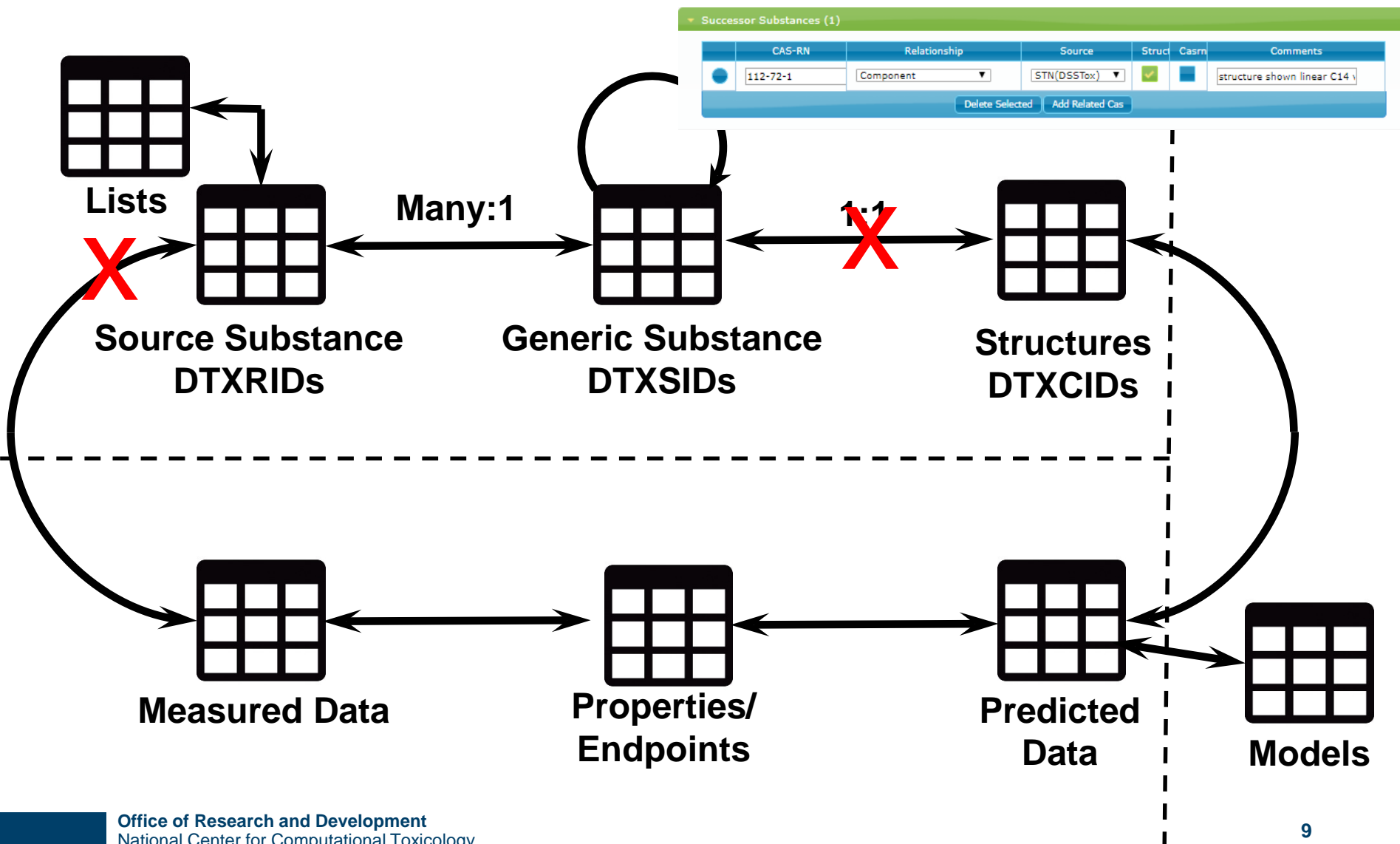
### Chemical Substances of Unknown or Variable Composition, Complex Reaction Products and Biological Materials (UVCB Substance) on the TSCA Inventory

This paper is a compendium of information related to the broad class of chemical substances referred to as UVCBs for the Toxic Substances Control Act (TSCA) Chemical Substance Inventory. These chemical substances cannot be represented by unique structures and molecular formulas.

- UVCB chemical examples
  - Surfactants with undefined composition
  - Petroleum Distillates
  - Gelatins, hydrozylates
  - Formaldehyde, reaction products with diethanolamine
  - Fatty acids, linseed-oil, compds. with triethylamine



# Data linkage breakdown



# So Where is the Data?

## C10-16 Alcohols

67762-41-8 | DTXSID4028331

**i** Searched by Approved Name: Found 1 result for 'C10-16 Alcohols'.

Presence in Lists

Federal

Safer Choice Chemical List

US State

International

Other

TSCAACTIVE

Record Information

Quality Control Notes

Related Substances

Synonyms

No Chemical Properties Found.

Related Substances

Synonyms

Links

Bioassays

Exposure

Hazard

Comments

Chemical Properties

Literature

Download / Send

Sort by: Relationship



Searched Chemical

1 related chemical  
structure with this  
substance

C10-16 Alcohols  
67762-41-8

Representative Component



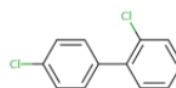
1-Tetradecanol  
112-72-1

# Chemical Families (e.g. PCBs)

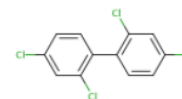
## ▼ Successor Substances (209)

	CAS-RN	Relationship
●	32774-16-6	is a Representative Isomer of this
●	2051-60-7	is a Representative Isomer of this
●	2051-61-8	is a Representative Isomer of this
●	2051-62-9	is a Representative Isomer of this
●	13029-08-8	is a Representative Isomer of this
●	16605-91-7	is a Representative Isomer of this
●	25569-80-6	is a Representative Isomer of this
●	33284-50-3	is a Representative Isomer of this
●	34883-43-7	is a Representative Isomer of this
●	34883-39-1	is a Representative Isomer of this
●	33146-45-1	is a Representative Isomer of this
●	2050-67-1	is a Representative Isomer of this
●	2974-92-7	is a Representative Isomer of this
●	2974-90-5	is a Representative Isomer of this
●	34883-41-5	is a Representative Isomer of this
●	2050-68-2	is a Representative Isomer of this

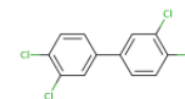
Download as: TSV Excel SDF



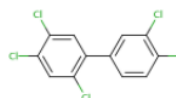
2,4-Dichlorobiphenyl  
34883-43-7



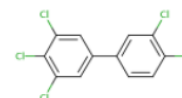
2,2',4,4'-Tetrachlorobiphenyl  
2437-79-8



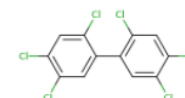
3,3',4,4'-Tetrachlorobiphenyl  
32508-13-3



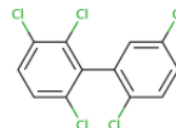
2,3',4,4',5-Pentachlorobiphenyl  
31508-00-8



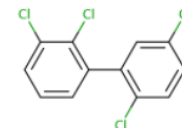
3,3',4,4',5-Pentachlorobiphenyl  
57465-28-8



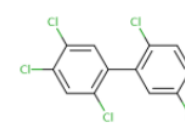
2,2',4,4',5,5'-Hexachlorobiphenyl  
35085-27-1



2,2',3,5',6-Pentachlorobiphenyl  
38379-99-8



2,2',3,5'-tetrachlorobiphenyl  
41464-39-5



2,2,4,5,5'-Pentachlorobiphenyl  
37680-73-2

Public ▼

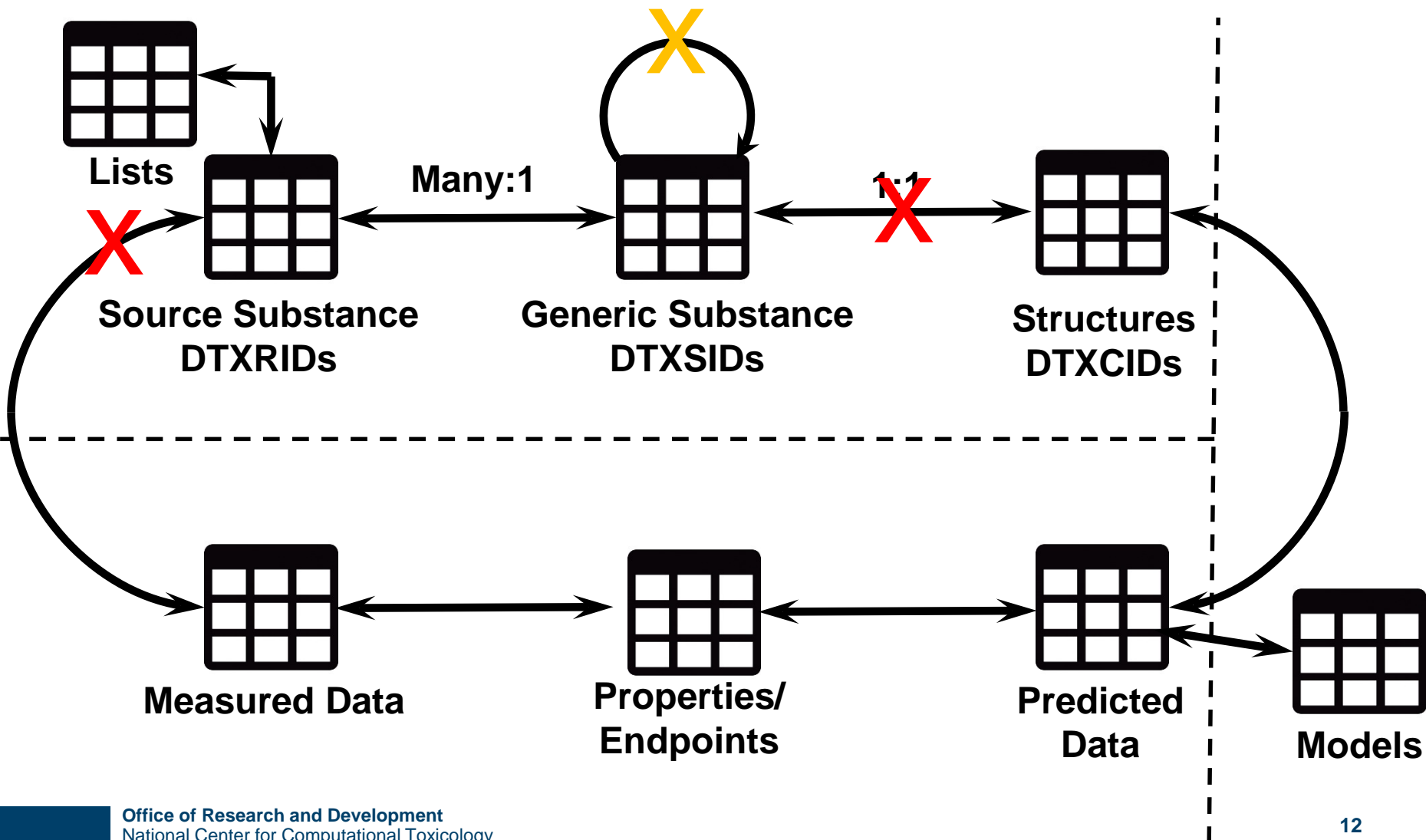
Public ▼

Public ▼

## Related Chemicals

Found 209 chemicals

# Data linkage breakdown





 tony27587@gmail.com 

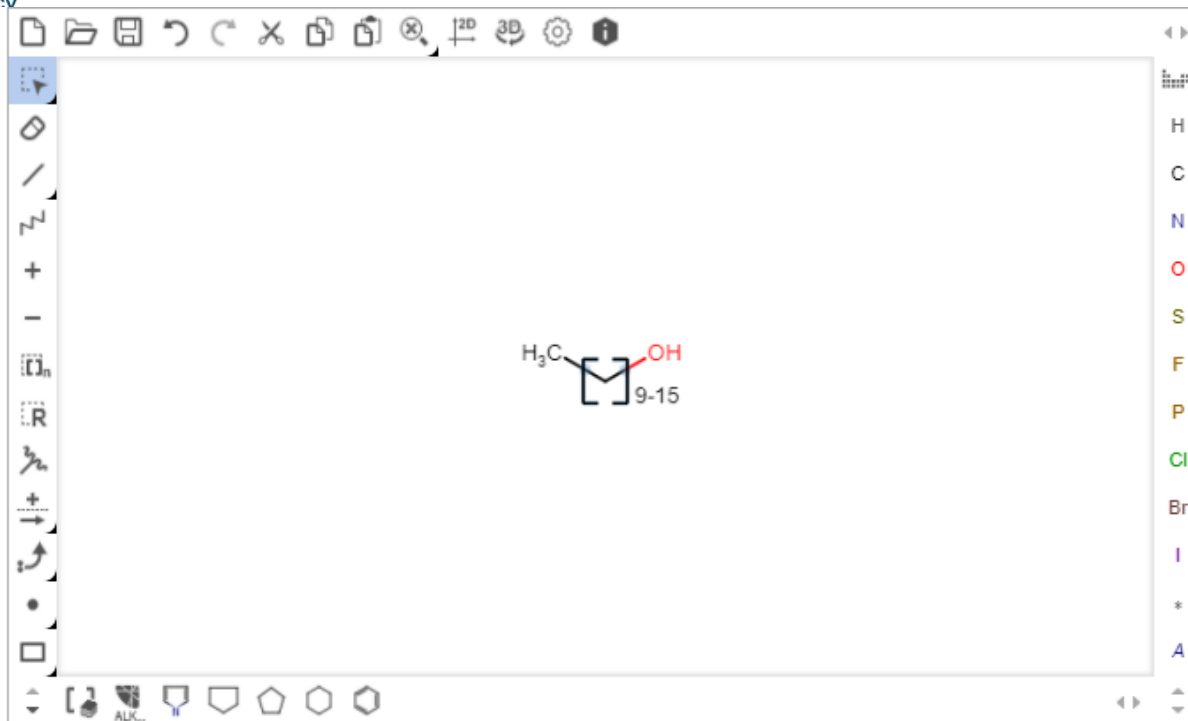


Markush Technology

[Download the Editor](#)



# Update Record



## Calculate from Structure

Substance\_ID: DTXSID4028331

CAS: 67762-41-8

Name: C10-16 Alcohols

Substance Type: Mixture/Formulation

QC Level: DSSTox\_High

Data Source: STN(DSSTox)

QC Notes:

SDA (Soap and Detergent Association) Reporting Number: 15-060-00. SDA Substance Name: C10-C16 alkyl alcohol

Compound\_ID:

Chemical Shown:

Markush Enumerable

Private Notes:

Source of CAS-Compound:

Double Stereo:

Chiral Stereo:

Chemical Form:

Public

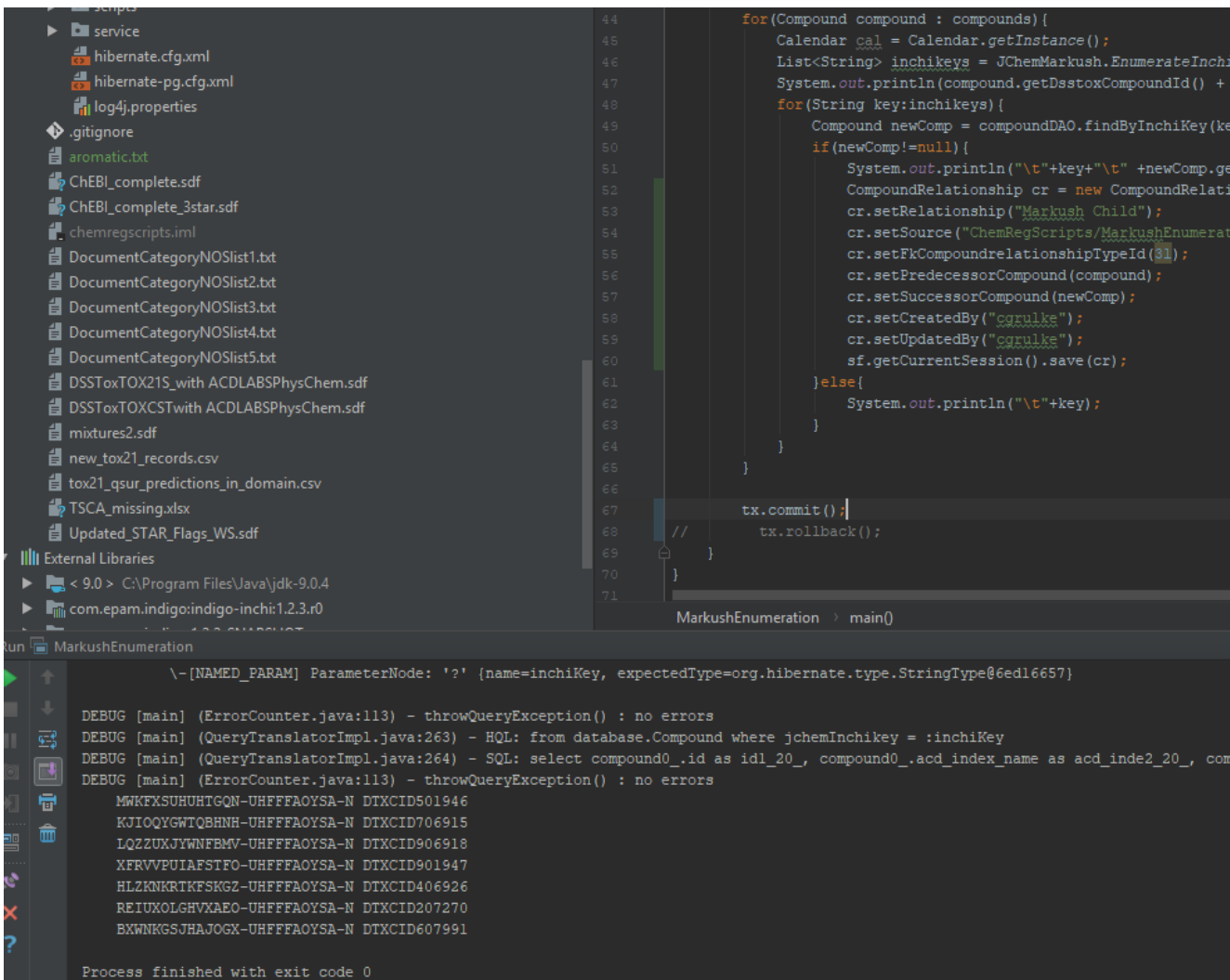
None

None

Organic

Organic Form: Parent

# Look Away: Work In Progress



```

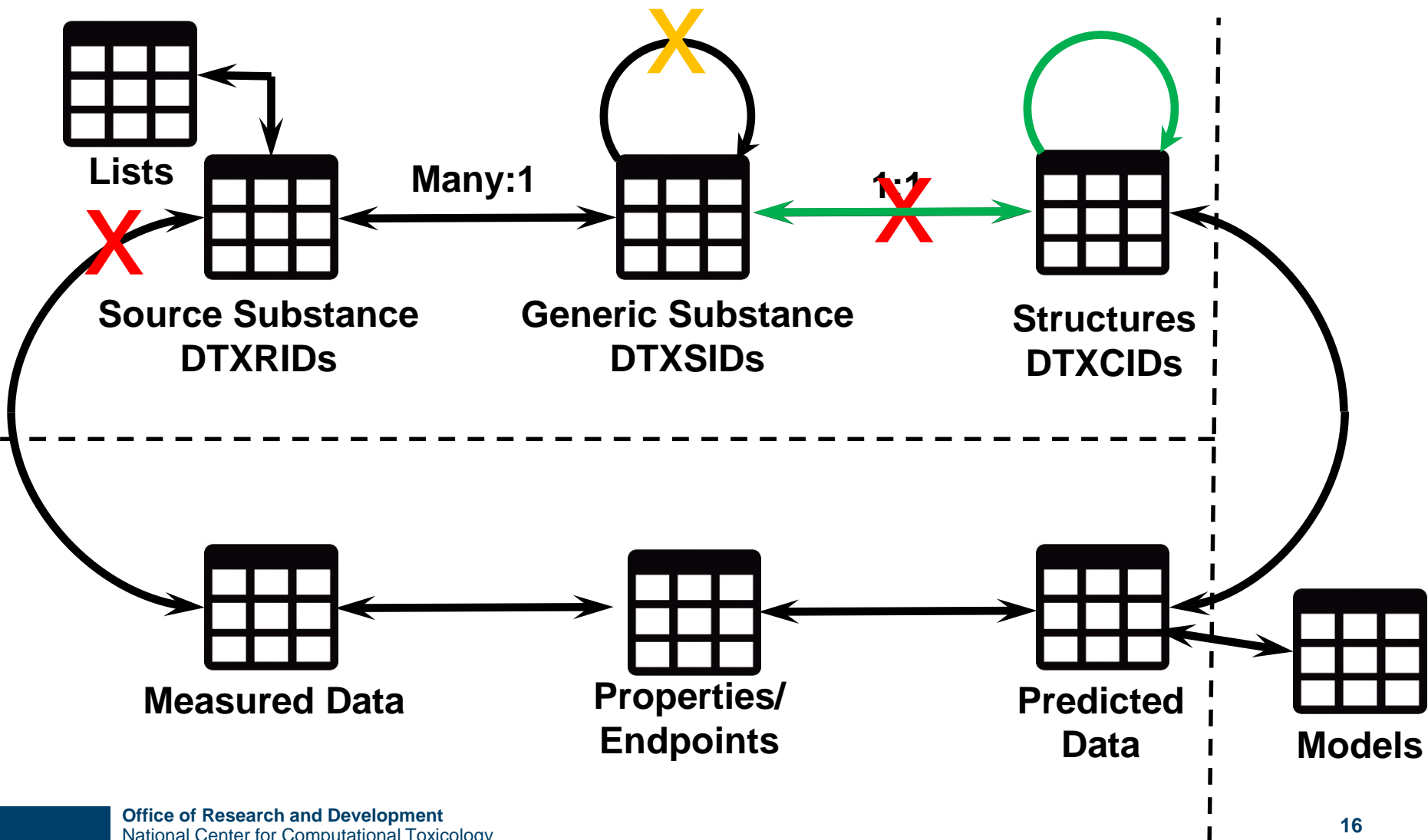
44     for(Compound compound : compounds){
45         Calendar cal = Calendar.getInstance();
46         List<String> inchikeys = JChemMarkush.EnumerateInchi
47         System.out.println(compound.getDsstoxCompoundId() +
48         for(String key:inchikeys){
49             Compound newComp = compoundDAO.findByInchiKey(ke
50             if(newComp!=null){
51                 System.out.println("\t"+key+"\t" +newComp.ge
52                 CompoundRelationship cr = new CompoundRelati
53                 cr.setRelationship("Markush Child");
54                 cr.setSource("ChemRegScripts/MarkushEnumerat
55                 cr.setFkCompoundrelationshipTypeId(31);
56                 cr.setPredecessorCompound(compound);
57                 cr.setSuccessorCompound(newComp);
58                 cr.setCreatedBy("cgrulke");
59                 cr.setUpdatedBy("cgrulke");
60                 sf.getCurrentSession().save(cr);
61             }else{
62                 System.out.println("\t"+key);
63             }
64         }
65     }
66     tx.commit();
67     // tx.rollback();
68 }
69 }
70 }
71 }
MarkushEnumeration > main()

\-[NAMED_PARAM] ParameterNode: '?' {name=inchiKey, expectedType=org.hibernate.type.StringType@6ed16657}

DEBUG [main] (ErrorCounter.java:113) - throwQueryException() : no errors
DEBUG [main] (QueryTranslatorImpl.java:263) - HQL: from database.Compound where jchemInchikey = :inchiKey
DEBUG [main] (QueryTranslatorImpl.java:264) - SQL: select compound0_.id as id1_20_, compound0_.acd_index_name as acd_inde2_20_, com
DEBUG [main] (ErrorCounter.java:113) - throwQueryException() : no errors

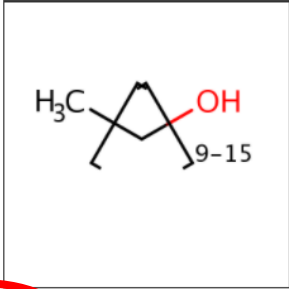
MWKFXSUHHTGQN-UHFFFAOYSA-N DTXCID501946
KJIOQYGTQBHNH-UHFFFAOYSA-N DTXCID706915
LQZZUXJYWNFBMV-UHFFFAOYSA-N DTXCID906918
XFRVVPUIAFSTFO-UHFFFAOYSA-N DTXCID901947
HLZKNKRIKFSKGZ-UHFFFAOYSA-N DTXCID406926
REIUXOLGHVXAEO-UHFFFAOYSA-N DTXCID207270
BXWNKGSJHAJOGX-UHFFFAOYSA-N DTXCID607991

Process finished with exit code 0
    
```





# New Relationships in Comptox Chemistry Dashboard



Related Substances | **Chemical Properties** | Analytical | Comments

**Intrinsic Properties**

Molecular Formula: (CH<sub>2</sub>)<sub>8</sub>-15CH<sub>3</sub>O Q Find

Average Mass: Not Found

Monoisotopic Mass: Not Found

**Structural Identifiers**

**Presence in Lists**

**Record Information**

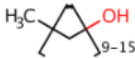
**Quality Control Notes**

Download / Send
Sort by: Relationship

8 chemicals


Hide:

Searched Chemical




C10-16 Alcohols  
67762-41-8

Markush Child




1-Decanol  
112-30-1

Markush Child




1-Tridecanol  
112-70-9

Markush Child




1-Undecanol  
112-42-6

Markush Child




1-Dodecanol  
112-53-8

Markush Child




1-Tetradecanol  
112-72-1

Markush Child



1-Pentadecanol  
629-78-5

Markush Child



1-Hexadecanol  
36653-82-4

# Collecting Property Distributions for Markush Children

Melting Point

Boiling Point

Surface Tension

Vapor Pressure

LogKoa: Octanol-Air

Henry's Law

Index of Refraction

Molar Refractivity

Molar Volume

Polarizability

## Experimental

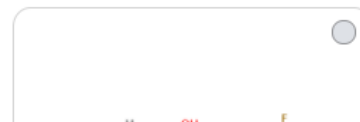
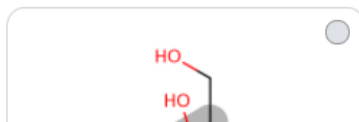
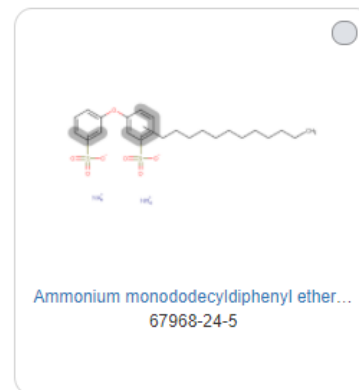
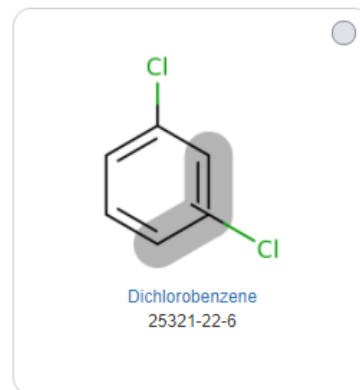
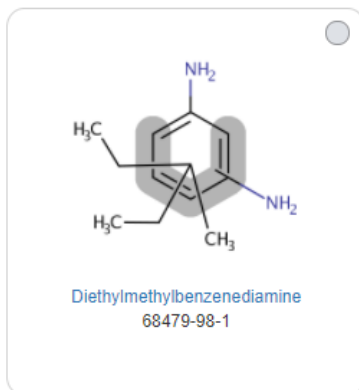
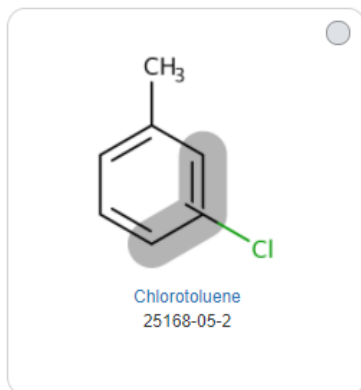
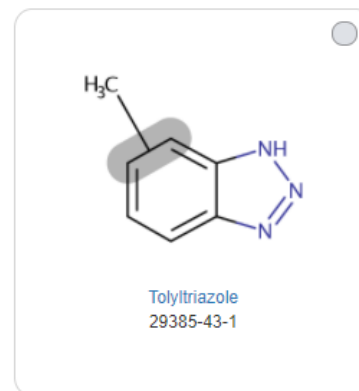
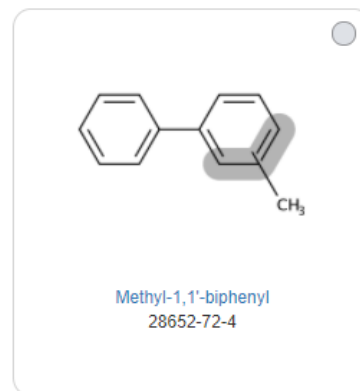
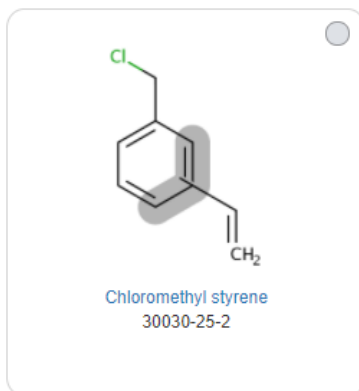
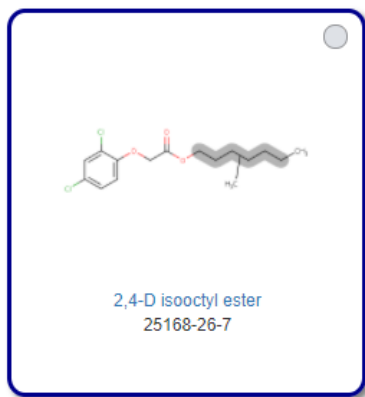
Source	Result	
PhysPropNCCT	2.52e-05 mmHg	

## Predicted

Source	Result	Calculated
NICEATM	5.23e-05 mmHg	Not Available
ACD/Labs	1.24e-04 mmHg	Not Available
TEST	1.24e-05 mmHg	TEST Report
OPERA	1.33e-05 mmHg	OPERA Model Report

# Markush in DSSTox

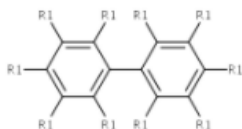
Download / Send Sort by: DTXSID 126 chemicals Hide: Select all



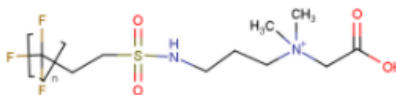
# Span of Predicted Properties for UVCBs

# Problems: Queryable, but Not Enumerable

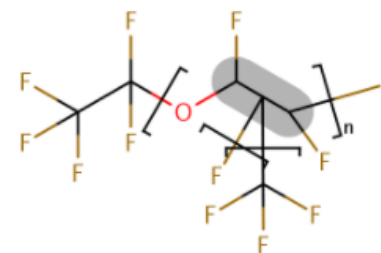
$R1 > 0, \text{restH}$



Polychlorinated biphenyls  
1336-36-3



Fluorotelomer Sulfonamido Betaines  
NOCAS\_892972



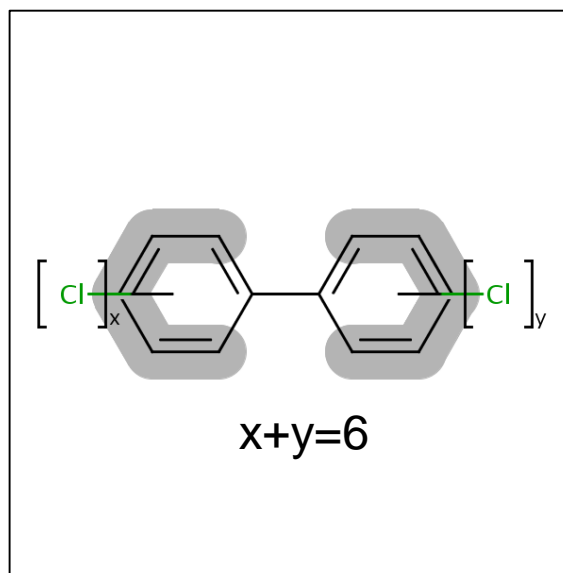
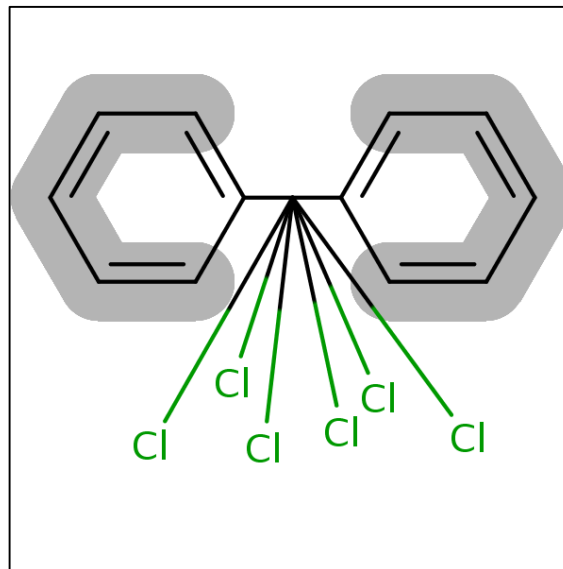
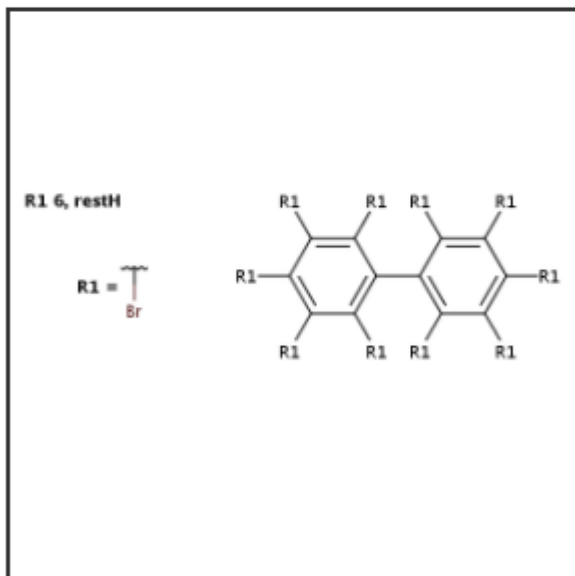
Krytox 143  
60164-51-4

# Problem: Human Readable or Computable

## Hexabromobiphenyl

36355-01-8 | DTXSID3025382

© Searched by Approved Name: Found 1 result for 'Hexabromobiphenyl'.



# Problem: Markush Uniqueness

Hexabrominated dibenzofurans

3635 NOCAS\_23924 | DTXSID3023924

© Searched by DSSTox\_Substance\_Id: Found 1 result for 'DTXSID3023924'.

R1 > 0, restH

R1 =

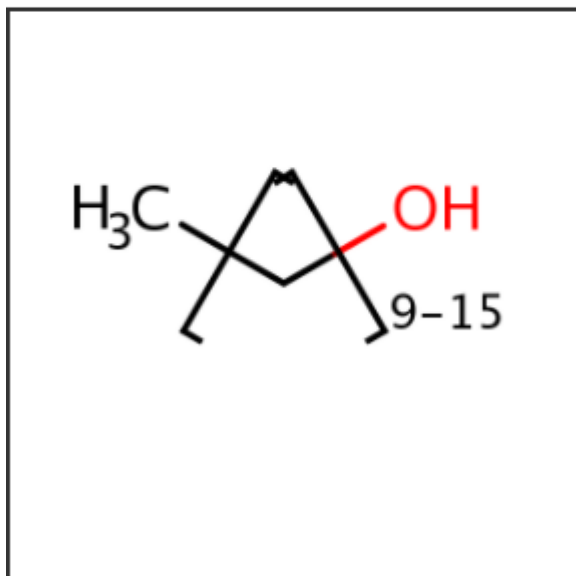
- InChI-Keys for Markush?
- Canonical Extend Smiles?
- Enumerate and compare children?

# Problem: Do We Know the Substance?

## C10-16 Alcohols

67762-41-8 | DTXSID4028331

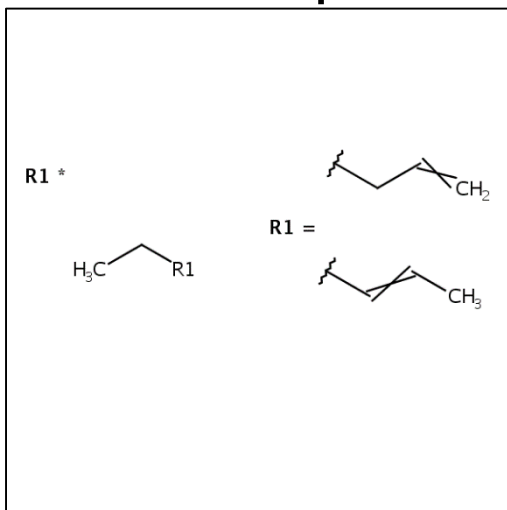
🔍 Searched by Approved Name: Found 1 result for 'C10-16 Alcohols'.





# Problem: Drawing Limitations

- Variable position bonds (e.g, Pentene)



- Mol File

- CXSMILES

ChemistryDashboard-Batch-Search\_2018-03-16\_13\_19\_53.sdf - MarvinView 18.5

File Edit View Table Structure Tools Help

#	structure	INPUT	FOUND_BY	DTXSID	PREFERRED_NAME
1		DTXSID40891095	DSSTox_Substance_id	DTXSID40891095	Benzenesulfonamide, 4-[[[4-[2-[4-(2-hydroxybutoxy)methyl]phenoxy]phenyl]amino]phenyl]sulfonamide sodium salt (1.1)
2		DTXSID90890210	DSSTox_Substance_id	DTXSID90890210	Benzenesulfonic acid, 3-[[[4-amino-9,10-dihydro-9,10-dioxo-2H-benzimidazole-2-yl]methyl]sulfonic acid sodium salt (1.2)
3		DTXS			Xylenes
4		DTXS			2-Mercaptomethylbenzimidazole
5		DTXSID4052184	DSSTox_Substance_id	DTXSID4052184	Glycerol dicaprate

Error

Read error after molecule 126:  
Invalid counts line

OK Stack Trace Copy to Clipboard

**Office of Research and Development**  
National Center for Computational Toxicology

# Conclusions

- Use of Markush for UVCB enables linking
- A lot of UVCB cannot be depicted with Markush
- Integration of Markush Technology brings its own problems
  - Best depiction vs current enumeration capabilities
  - Determining uniqueness
  - Actually understanding the chemical substance
  - Overcoming limitations through willpower
  - Sharing Markush

# Acknowledgements



Credit: the Research Triangle Foundation

EPA NCCT IT  
Jeff Edwards  
Jeremy Dunne

EPA NCCT Curation  
David McKee  
Inthirany Thillainadarajah  
Sakuntala Sivasupramaniam

EPA NERL Curation  
Brian Meyer

# Questions?





# Conclusion

- The CompTox Chemistry Dashboard is the public interface into DSSTox chemical data (and everything that we can associate)
- Semi-automated curation is effective for expanding our data, while allowing appropriate correction via manual interaction
- We take care of our data
- We let you know how well we take care of our data
- We want your feedback to help us take care of our data!