

Mixtures QSAR

modelling collections of chemicals

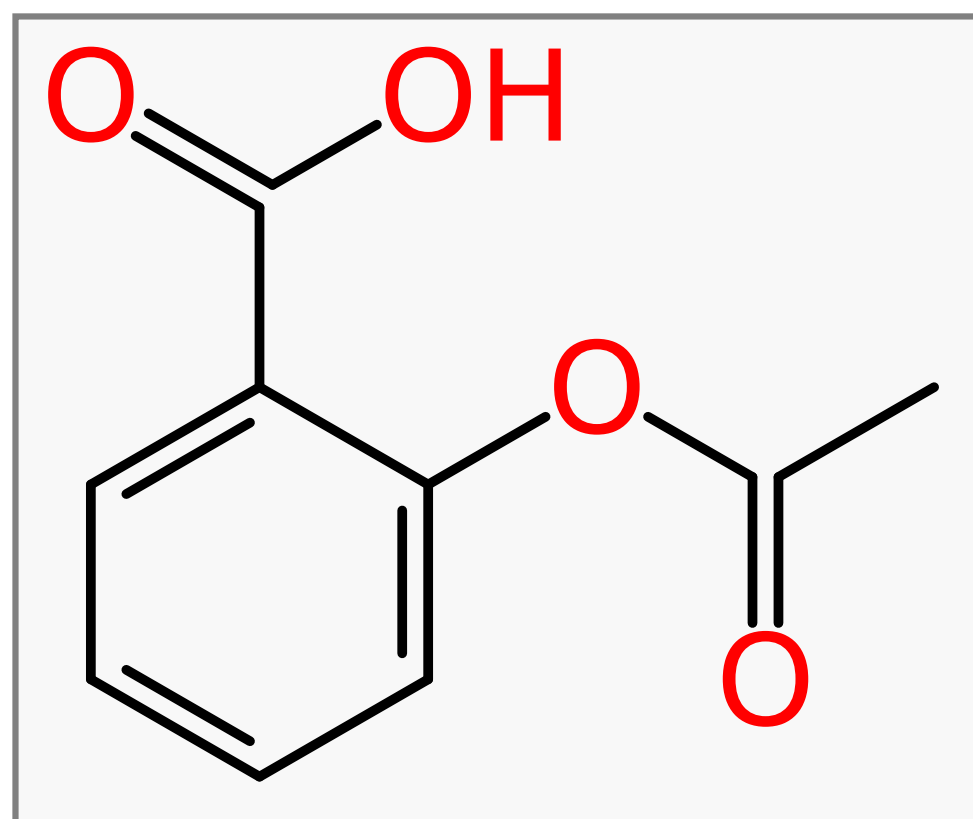
Alex M. Clark

alex@collaborativedrug.com



CDD.VAULT[®]
Complexity Simplified

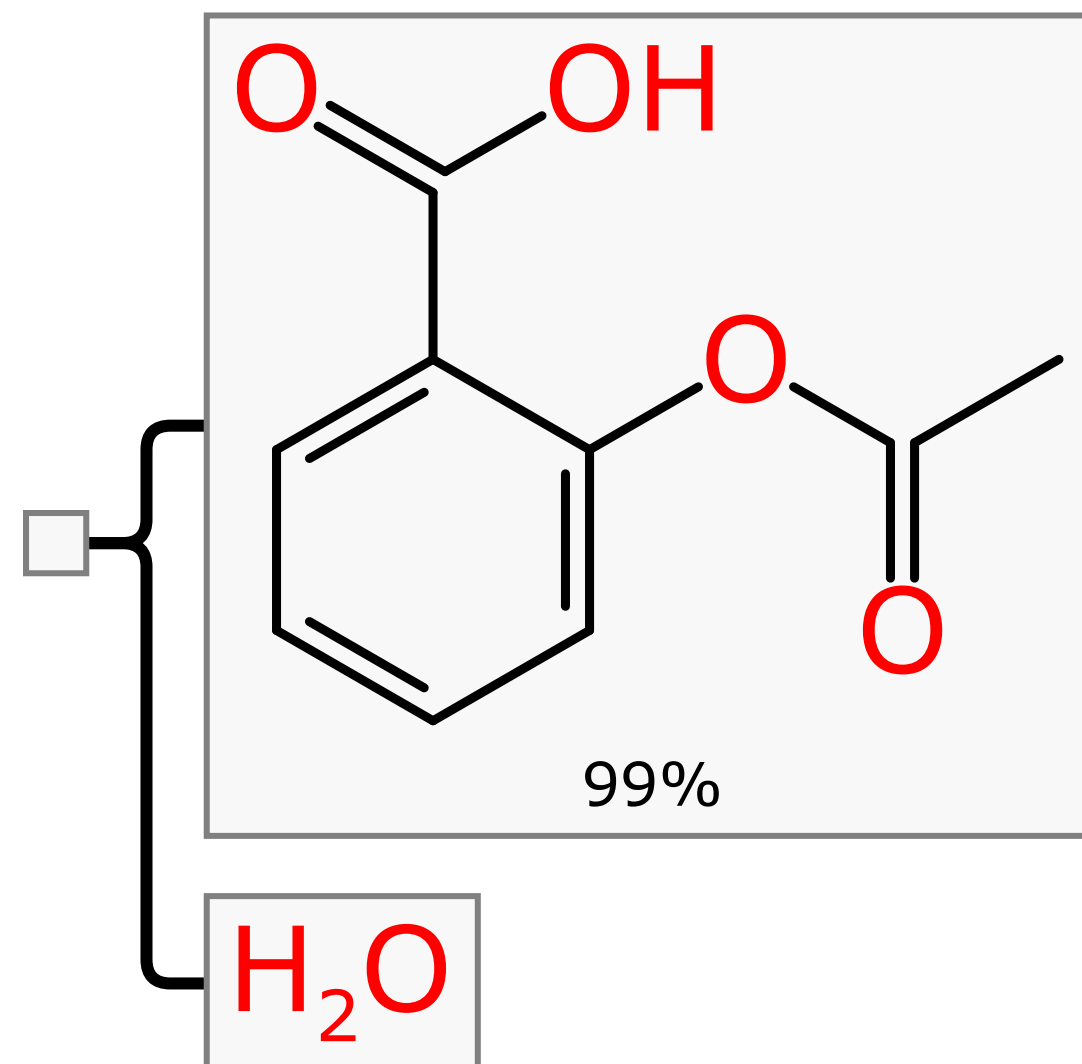
Real world chemistry is mixtures



Molfile

InChI

1980's



Mixfile

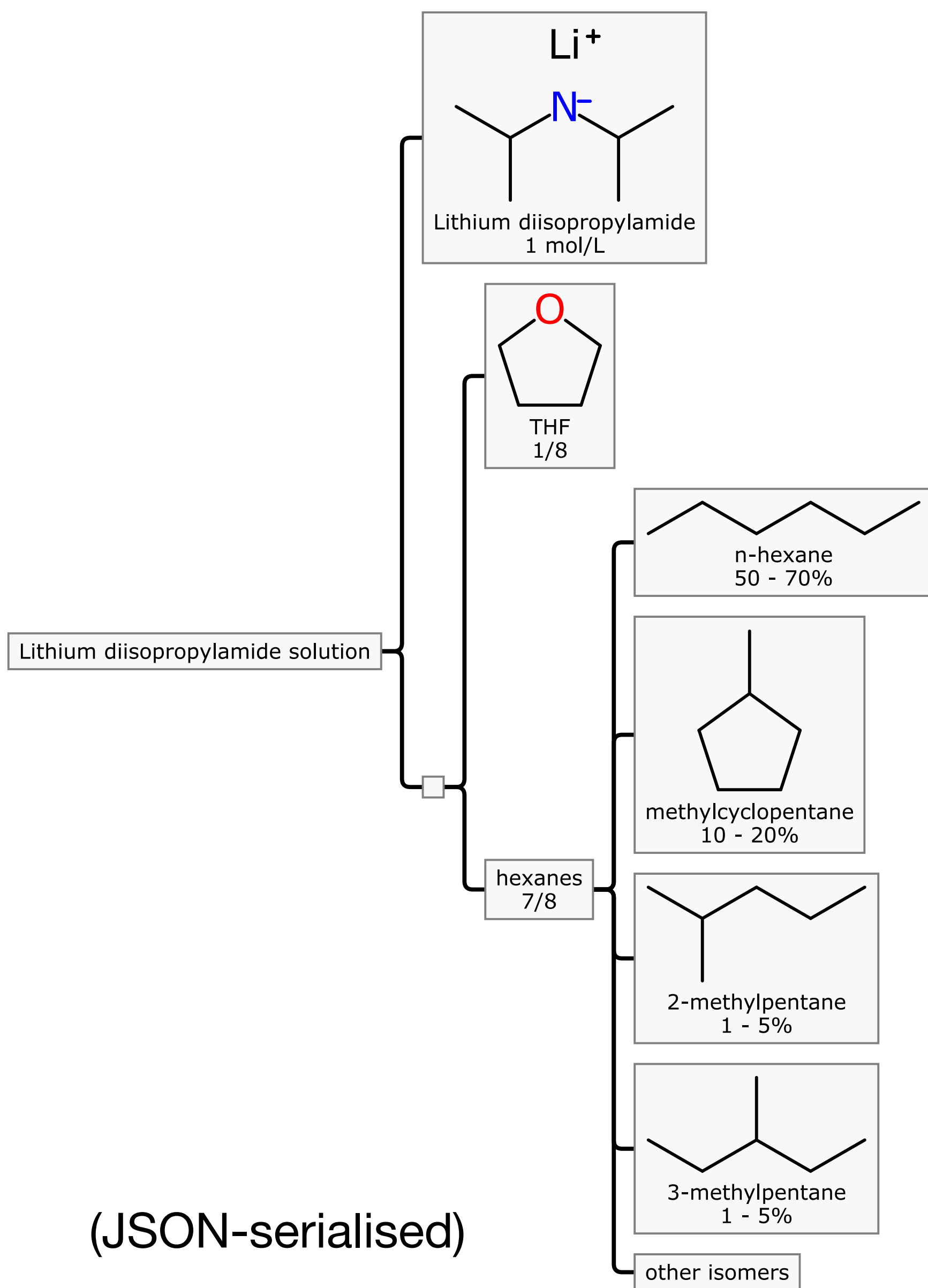
MInChI

2020's

- ❖ Most mixtures stored as text or custom spreadsheets
- ❖ Value of upgrading to cheminformatics is well established...
- ❖ ... with the right datastructure, always just one script away from what you need
- ❖ If you can represent it, you can model it



Mixfile/MInChI



❖ Format needs to be:

- ▶ hierarchical
- ▶ embed structures when possible
- ▶ include concentration information
- ▶ tolerate uncertainty

❖ More verbose ELN-friendly form is **Mixfile**

❖ Concise form with canonical components is **MInChI** (*mixtures* InChI) notation

```
MInChI=0.00.1S/C4H8O/c1-2-4-5-3-1/h1-4H2&C6H12/  
c1-6-4-2-3-5-6/h6H,2-5H2,1H3&C6H14/c1-3-5-6-4-2/  
h3-6H2,1-2H3&C6H14/c1-4-5-6(2)3/h6H,4-5H2,1-3H3&C6H14/  
c1-4-6(3)5-2/h6H,4-5H2,1-3H3&C6H14N.Li/c1-5(2)7-6(3)4;/  
h5-6H,1-4H3;/q-1;+1/n{6&{1&{3&2&4&5}}}/  
g{1mr0&{1vp0&{5:7pp1&1:2pp1&1:5pp0&1:5pp0}7vp0}}
```

Journal of Cheminformatics (2019)

[10.1186/s13321-019-0357-4](https://doi.org/10.1186/s13321-019-0357-4)

Content creation

The screenshot shows the CDD.VAULT software interface for creating a mixture. The main window displays a 'Mixture' editor with a toolbar. A central panel shows search results for 'hexanes', including various chemical structures and their corresponding names and identifiers. The results include:

- Hexanes $\geq 98.5\%$
- Hexanes $\geq 99\%$
- Hexanes $\geq 98\%$
- HEXANES UNII 1S70G88A39
- Sodium 1-hexanesulfonate 99%
- Sodium 1-hexanesulfonate monohydrate 99%
- N-ETHYLCYCLOHEXANESULFAMIC ACID UNII 1LK9V1E23U CASRN 69103-74-8 PubChemCID 20316519

A 'Submit' button is visible at the bottom right of the search panel.

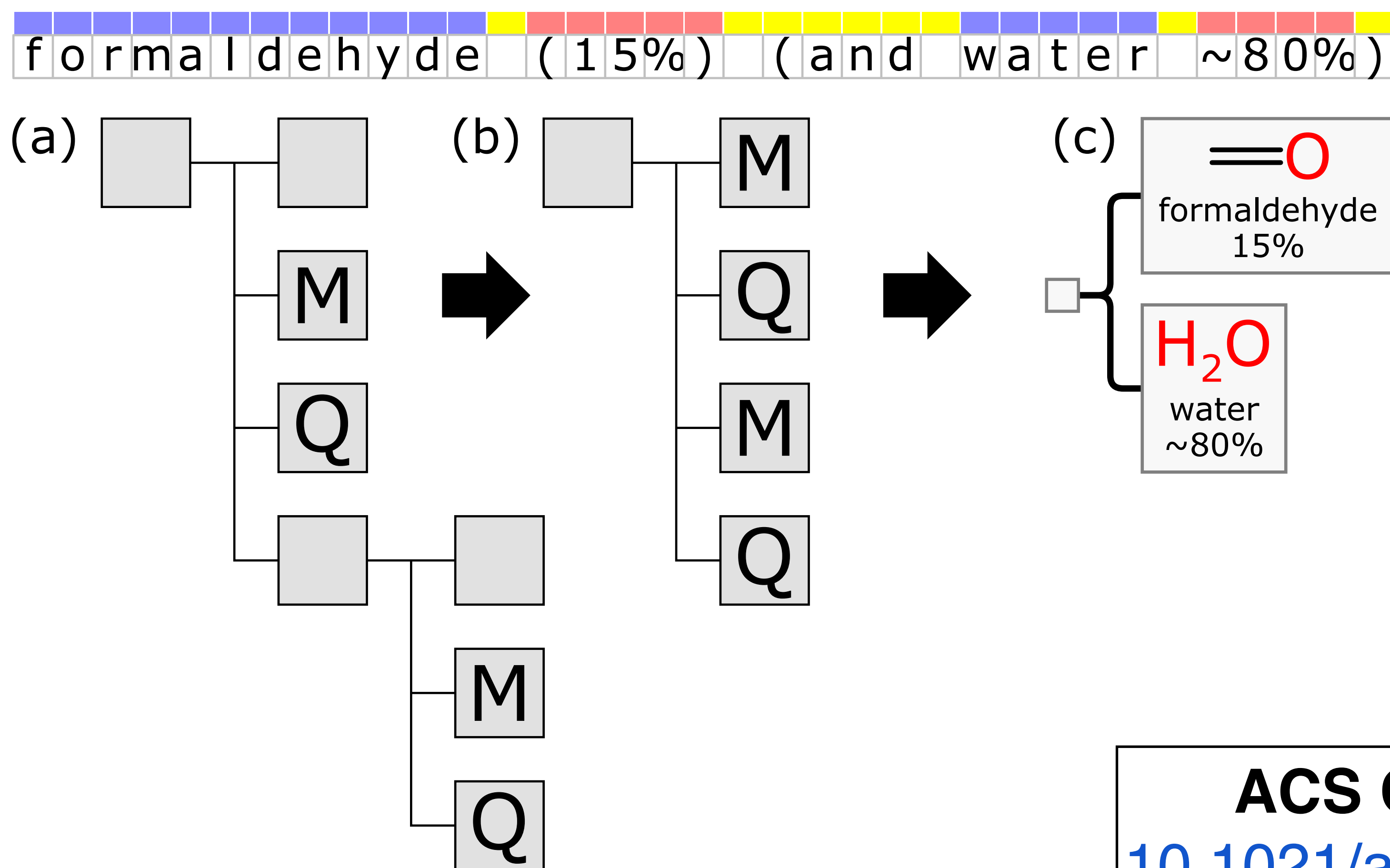
✿ Draw with editor

github.com/cdd/mixtures

✿ ... or import existing data

Fine chemicals as text

- ❖ Catalogs & inventories favour short descriptions of mixtures
- ❖ Machine learning for text-to-mixture works well:



Mixtures as spreadsheets

❖ Each component defined by a column (or several)...

The screenshot shows a Microsoft Excel spreadsheet with the following data:

	A	B	C	D	E	F	G	H	I	J	K	L
1	Code	Major	Conc (uM)	MPP-Cyto%	Minor1	Minor2	Minor3	Minor4	Minor5	Minor6	Minor7	Minor8
2	ENT1a	cannabidiol	10	62	cannabichromene	cannabigerol	cannabidivarin	limonene	linalool	nerolidol	alpha-pinene	cis-phytol
3	ENT2a	cannabidiol	10	16	limonene	linalool	nerolidol	alpha-pinene	cis-phytol			
4	ENT3a	cannabidiol	10	31	cannabichromene	cannabigerol	cannabidivarin					
5	ENT4a	cannabidiol	10	16	cannabichromene	limonene	linalool	nerolidol	alpha-pinene	cis-phytol		
6	ENT5a	cannabidiol	10	14	cannabigerol	limonene	linalool	nerolidol	alpha-pinene	cis-phytol		
7	ENT6a	cannabidiol	10	51	cannabidivarin	limonene	linalool	nerolidol	alpha-pinene	cis-phytol		
8	ENT7a	cannabidiol	10	20	cannabichromene	cannabigerol						
9	ENT8a	cannabidiol	10	44	cannabichromene	cannabidivarin						
10	ENT9a	cannabidiol	10	36	cannabigerol	cannabidivarin	limonene	linalool	nerolidol	alpha-pinene	cis-phytol	
11	ENT10a	cannabidiol	10	64	cannabichromene	cannabigerol	cannabidivarin	limonene	linalool			
12	ENT11a	cannabidiol	10	33	cannabichromene	cannabigerol	cannabidivarin	nerolidol				
13	ENT12a	cannabidiol	10	27	cannabichromene	cannabigerol	cannabidivarin	alpha-pinene	cis-phytol			
14	ENT1b	cannabinol	10	48	cannabichromene	cannabigerol	cannabidivarin	limonene	linalool	nerolidol	alpha-pinene	cis-phytol
15	ENT2b	cannabinol	10	6	limonene	linalool	nerolidol	alpha-pinene	cis-phytol			
16	ENT3b	cannabinol	10	41	cannabichromene	cannabigerol	cannabidivarin					
17	ENT4b	cannabinol	10	9	cannabichromene	limonene	linalool	nerolidol	alpha-pinene	cis-phytol		

CDD Vault

- ❖ Mixture composer
- ❖ Destination for each column
- ❖ Build mixture hierarchy from rows
- ❖ Name-to-structure, lookup databases, links to ID codes

CDD VAULT · McKerrow Vault

Help · Log out

Explore Data ELN **Import Data** Reports Settings 4 Full-Access User

Step 1: Choose Data File Step 2: Map Fields Step 3: Commit Data

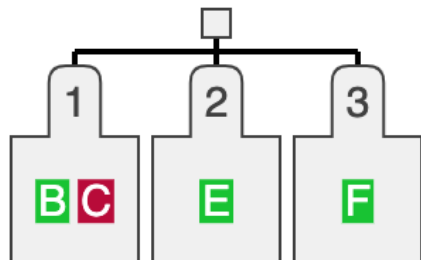
File: PhAROS_US20180098948A1.xlsx Project: Mixtures · Owner: Full-Access User

Compose mixtures from columns [Edit](#)

Not part of mixture
 Chemical structure [Molfile]
 Chemical structure [SMILES]
 Molecule lookup by name
 CDD Vault identifier
 Component name (text)
 Quantity
 Units

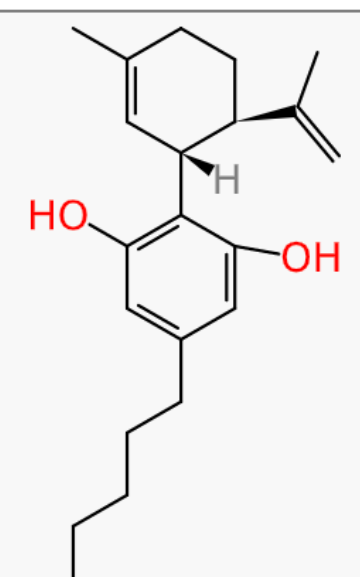
Component

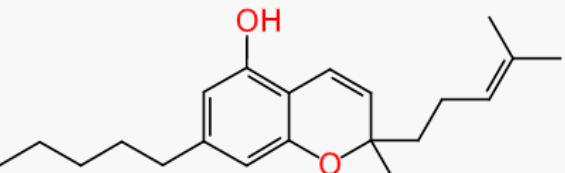
1	2	3
4	5	6
7	8	9

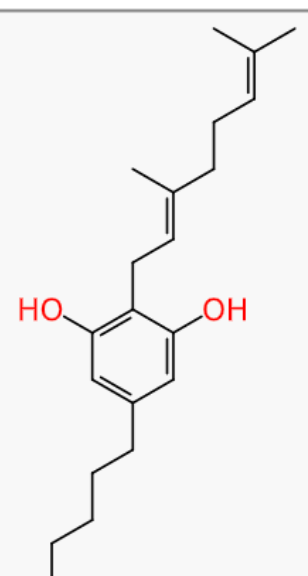


B ⇒ 1 molecule name	C ⇒ 1 quantity	D	E ⇒ 2 molecule name	F ⇒ 3 molecule name	G
Major	Conc (uM)	MPP-Cyto%	Minor1	Minor2	Minor3
cannabidiol	10	62	cannabichromene	cannabigerol	cannabidivari
cannabidiol	10	16	limonene	linalool	nerolidol
cannabidiol	10	31	cannabichromene	cannabigerol	cannabidivari
cannabidiol	10	16	cannabichromene	limonene	linalool
cannabidiol	10	14	cannabigerol	limonene	linalool
cannabidiol	10	51	cannabidivarin	limonene	linalool
cannabidiol	10	20	cannabichromene	cannabigerol	
cannabidiol	10	44	cannabichromene	cannabidivarin	
cannabidiol	10	36	cannabigerol	cannabidivarin	limonene
cannabidiol	10	64	cannabichromene	cannabigerol	cannabidivari

and 56 additional rows


Cannabidiol
10 µmol/L
UNII 19GBJ60SN5
CASRN 13956-29-1
PubChemCID 644019
COSING 96287
INCI CANNABIDIOL - DERIVED FROM EXTRAC...


CANNABICHROMENE
UNII K4497H250W
CASRN 20675-51-8
PubChemCID 30219


CANNABIGEROL
UNII J1K406072N
CASRN 25654-31-3
PubChemCID 5315659
COSING 98212
INCI CANNABIGEROL

High throughput screening

Mixture Export - pubchem504823.mixfile

Activity Protocol Results

1 2 3

1 pubchemAID 504823

2 60 µL

3 Hep-2 Cells substrate

4 H₂O water solvent

5 40 µL

6 virus substrate

7 H₂O water solvent

8 20 µL

9 SIM 0.39-50 µmol/L compound

10 DMSO 3% solvent

11 H₂O water solvent

Edit Template

Molecule structure 9

	1	2	3	4
PUBCHEM_RESULT_TAG	1	2	3	4
PUBCHEM_SID	104169543	104169547	118045990	118045991
PUBCHEM_CID	49842897	6619281	15989137	51051584
PUBCHEM_ACTIVITY_OUTCOME	Active	Active	Active	Active
PUBCHEM_ACTIVITY_SCORE	100	100	100	100
PUBCHEM_ACTIVITY_URL				
PUBCHEM_ASSAYDATA_COMMENT				
IC50 Modifier			<	
IC50	3.2	2.9	1.2	2.4
% CPE Inhibition @ 150 µM			70.12	6.22

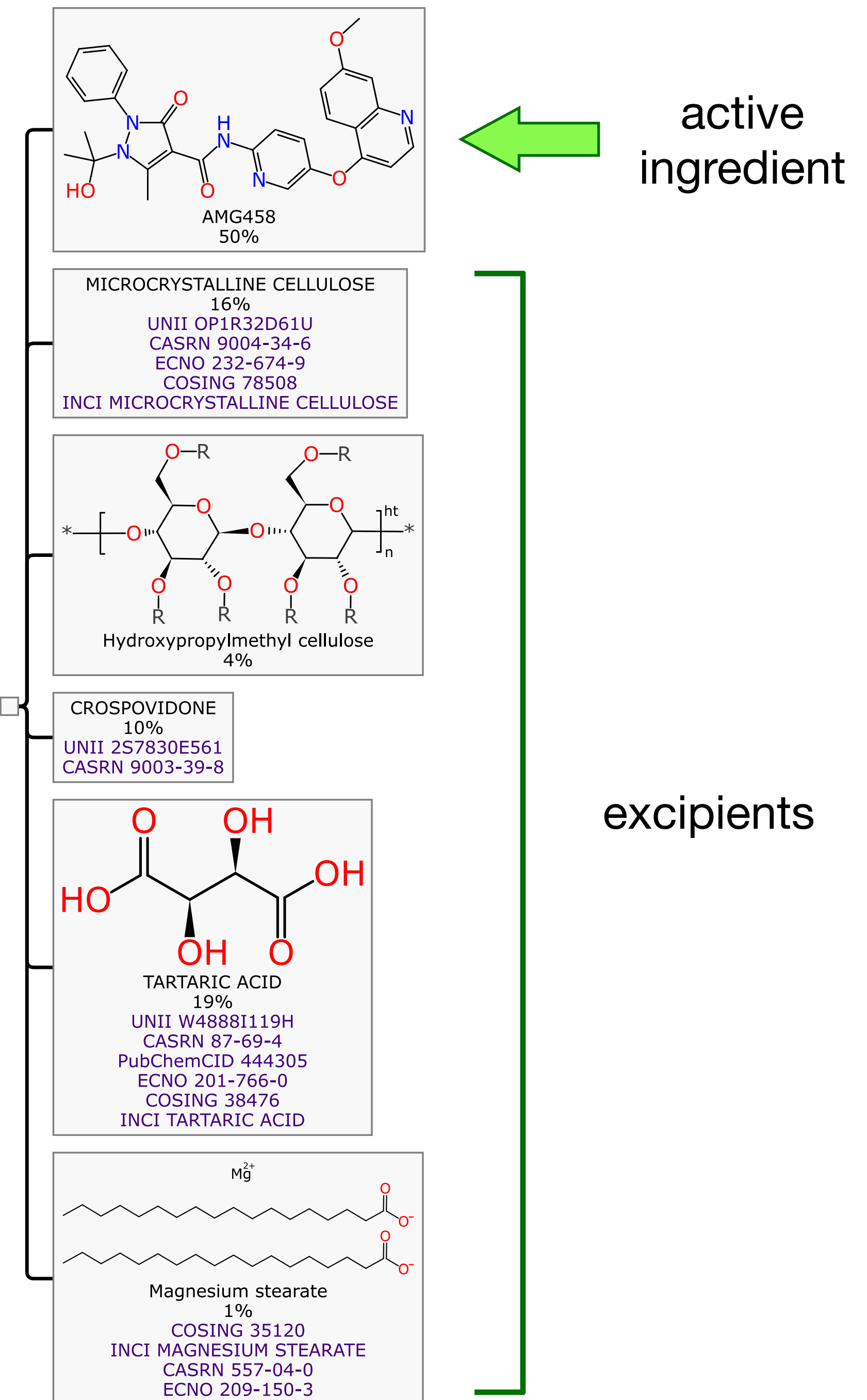
Drug tablet formulation

❖ Represent *active ingredient* by structure & concentration

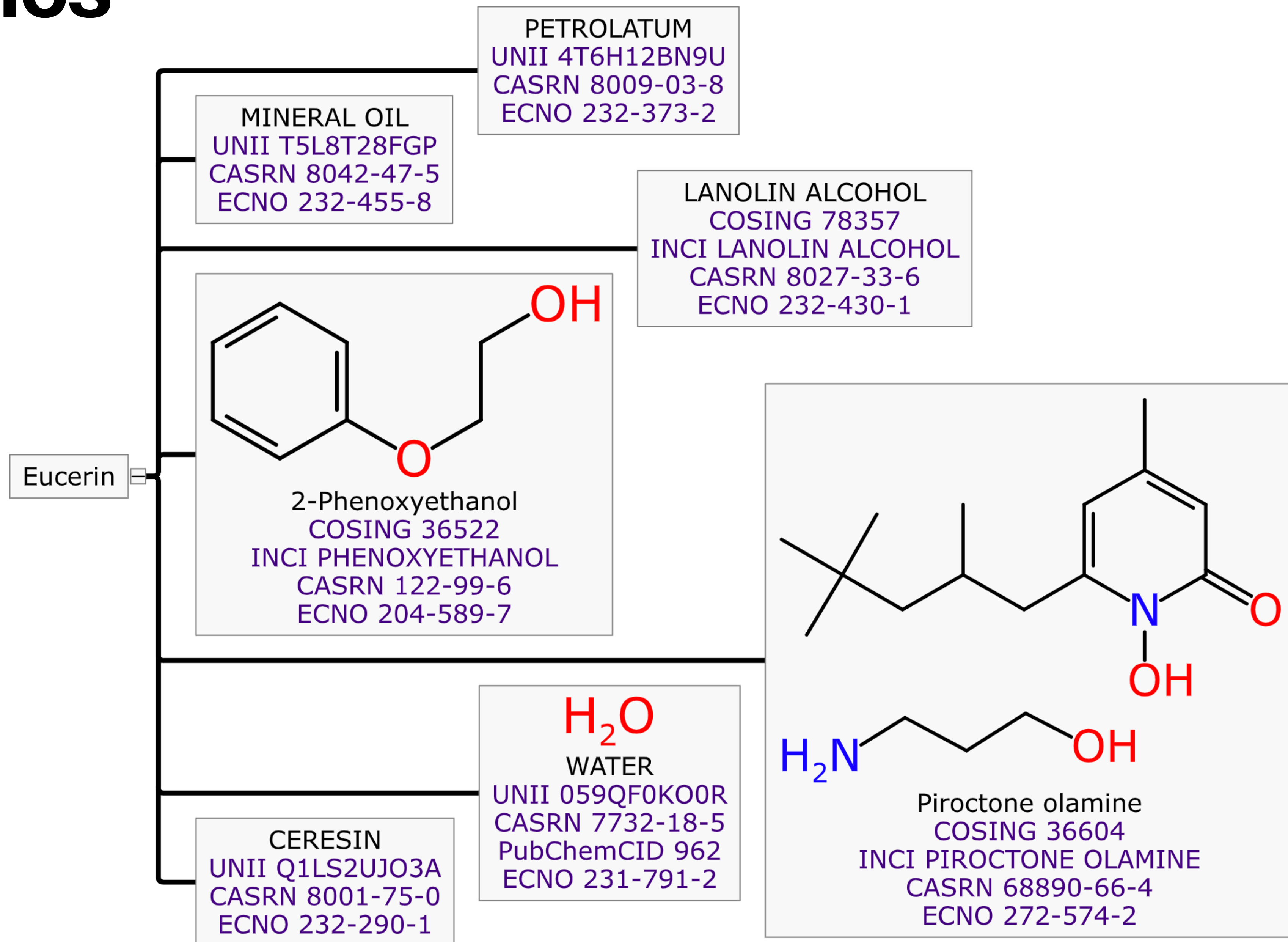
❖ Other materials (*excipients*):

- ▶ polymers
- ▶ organic molecules
- ▶ salts
- ▶ amorphous materials

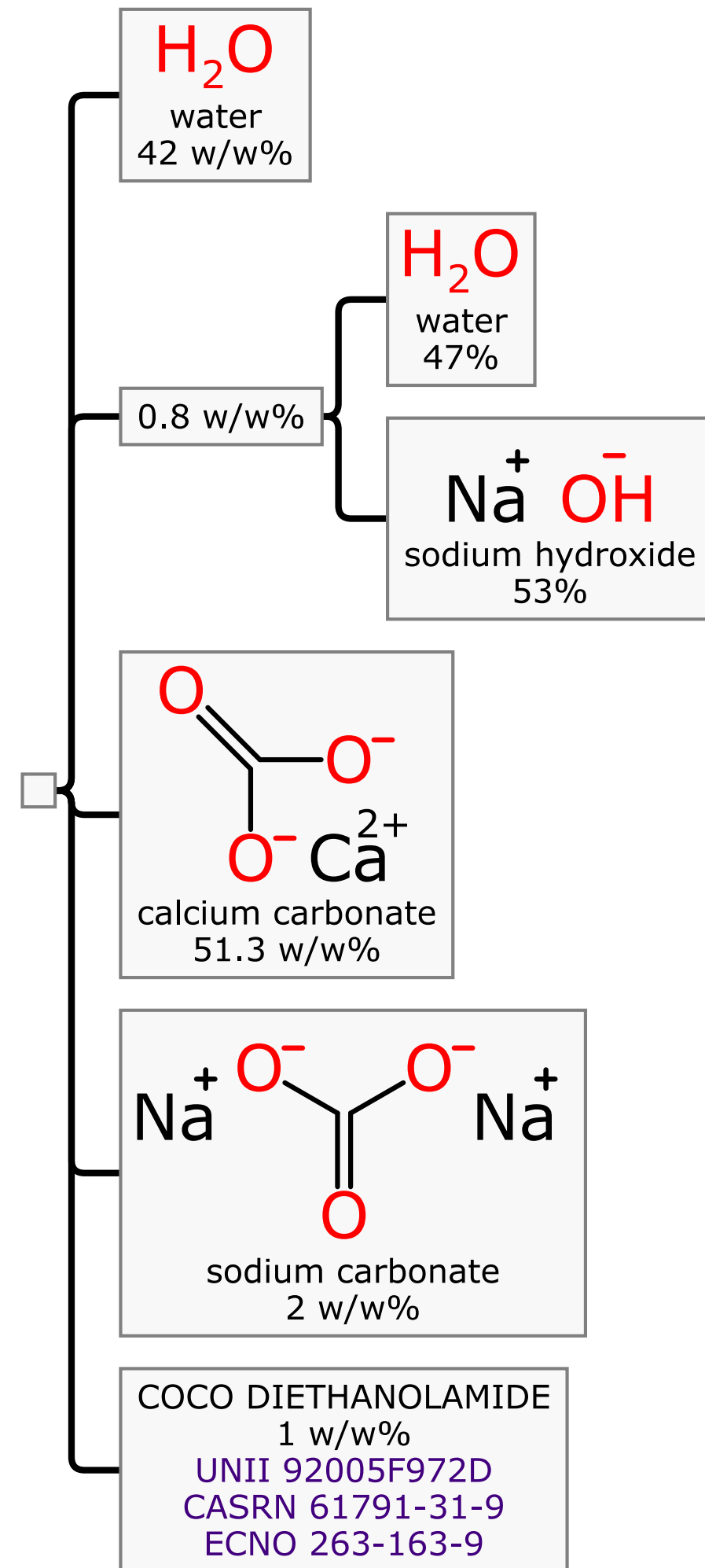
J Pharm Sci (2009)
[10.1002/jps.21422](https://doi.org/10.1002/jps.21422)



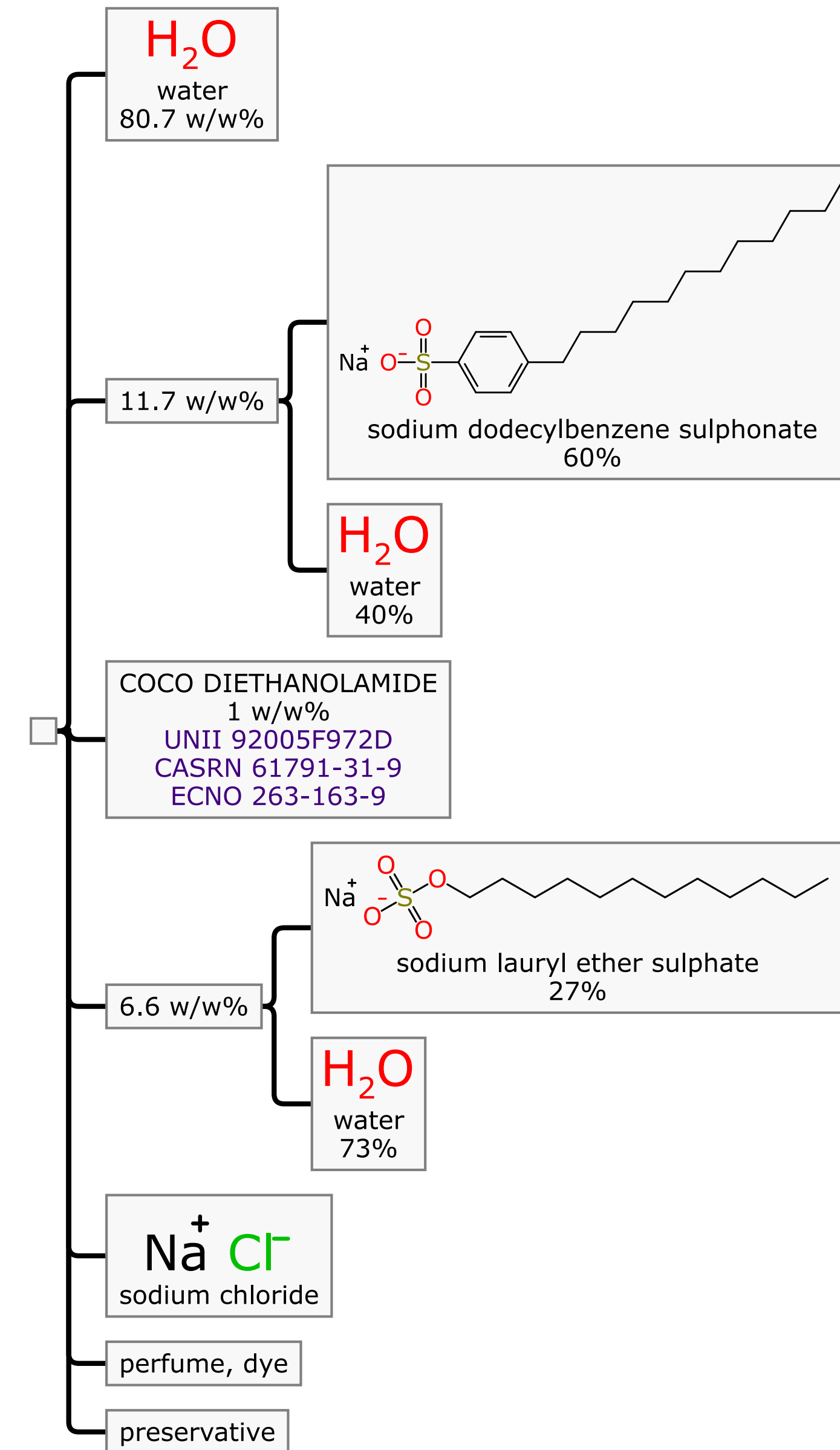
Cosmetics



Consumer Products

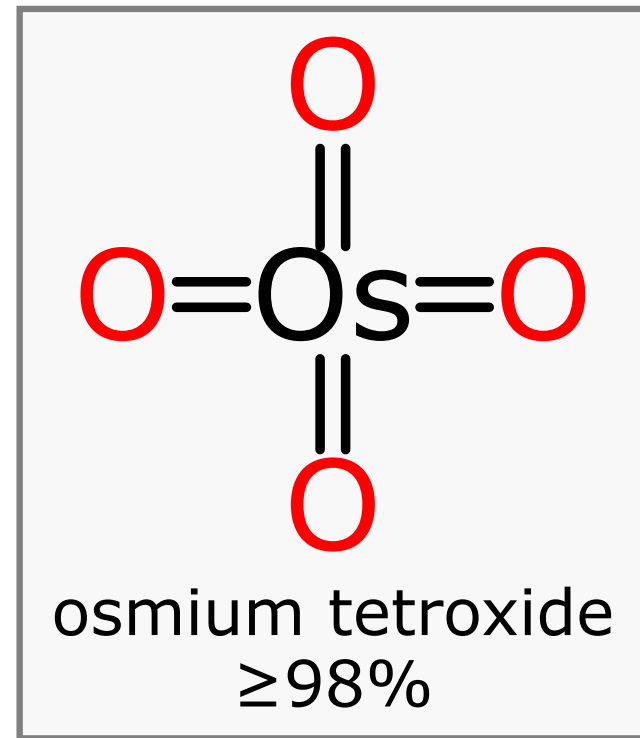


abrasive hand cleaner

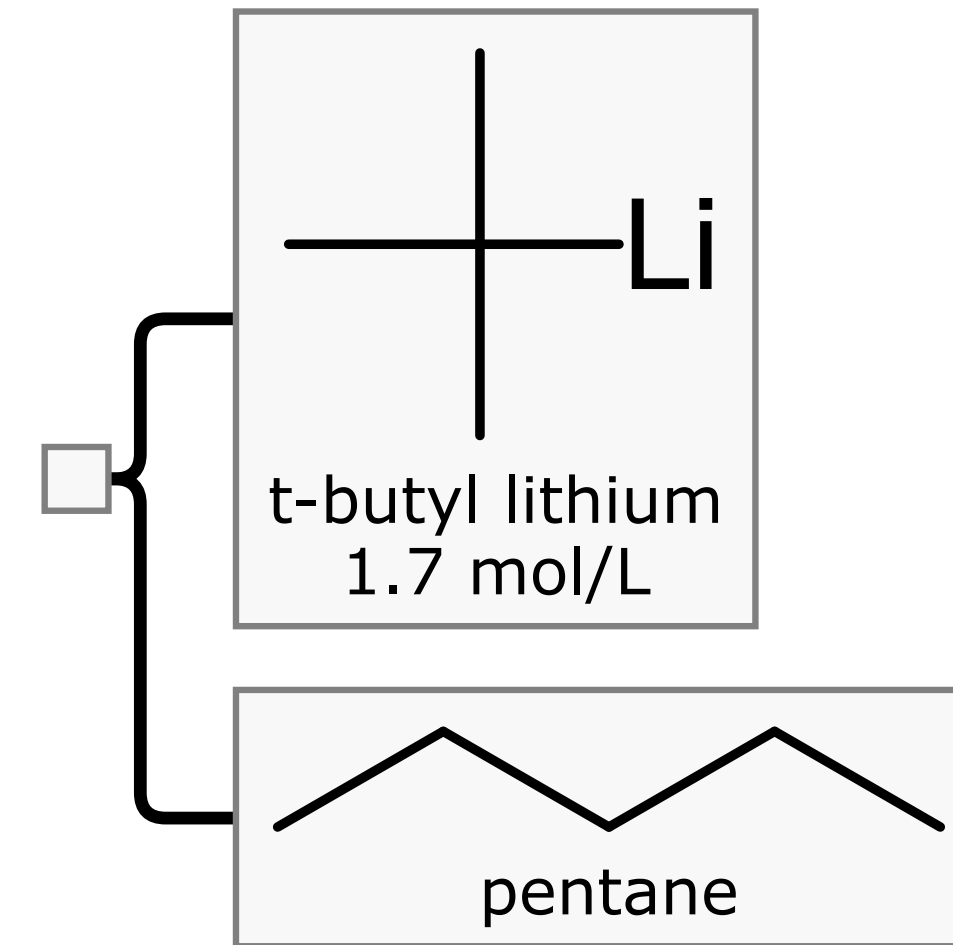


dishwashing liquid

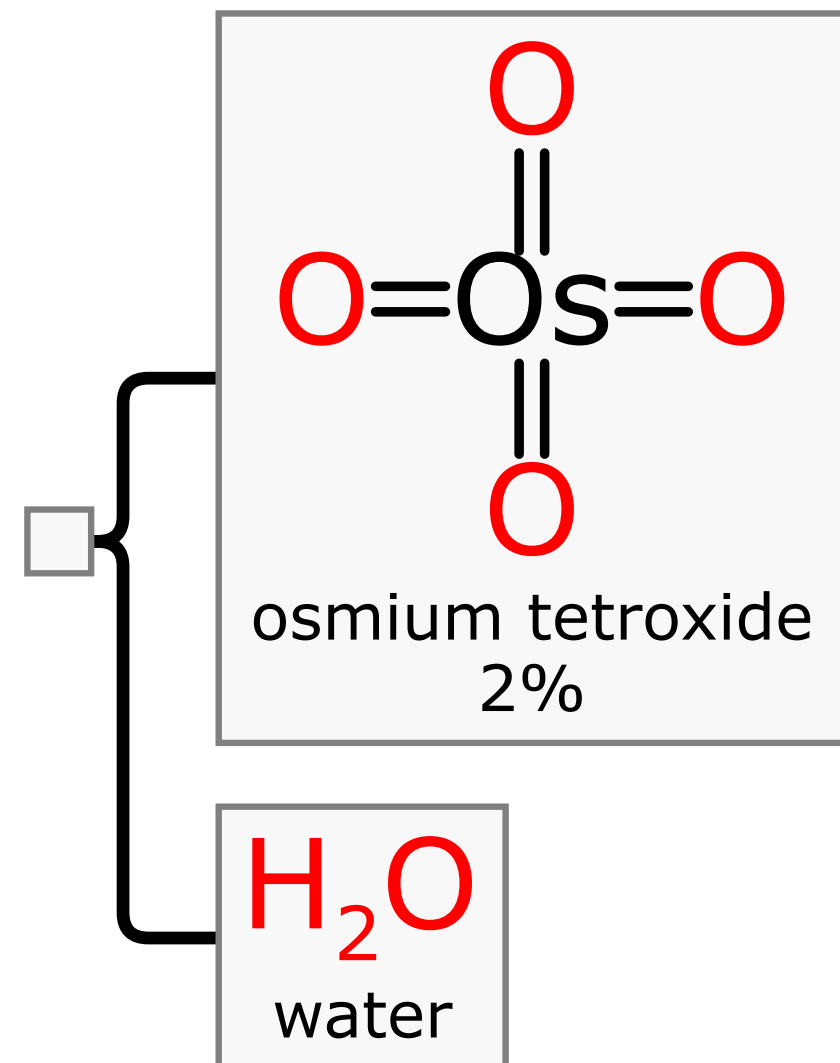
Hazards



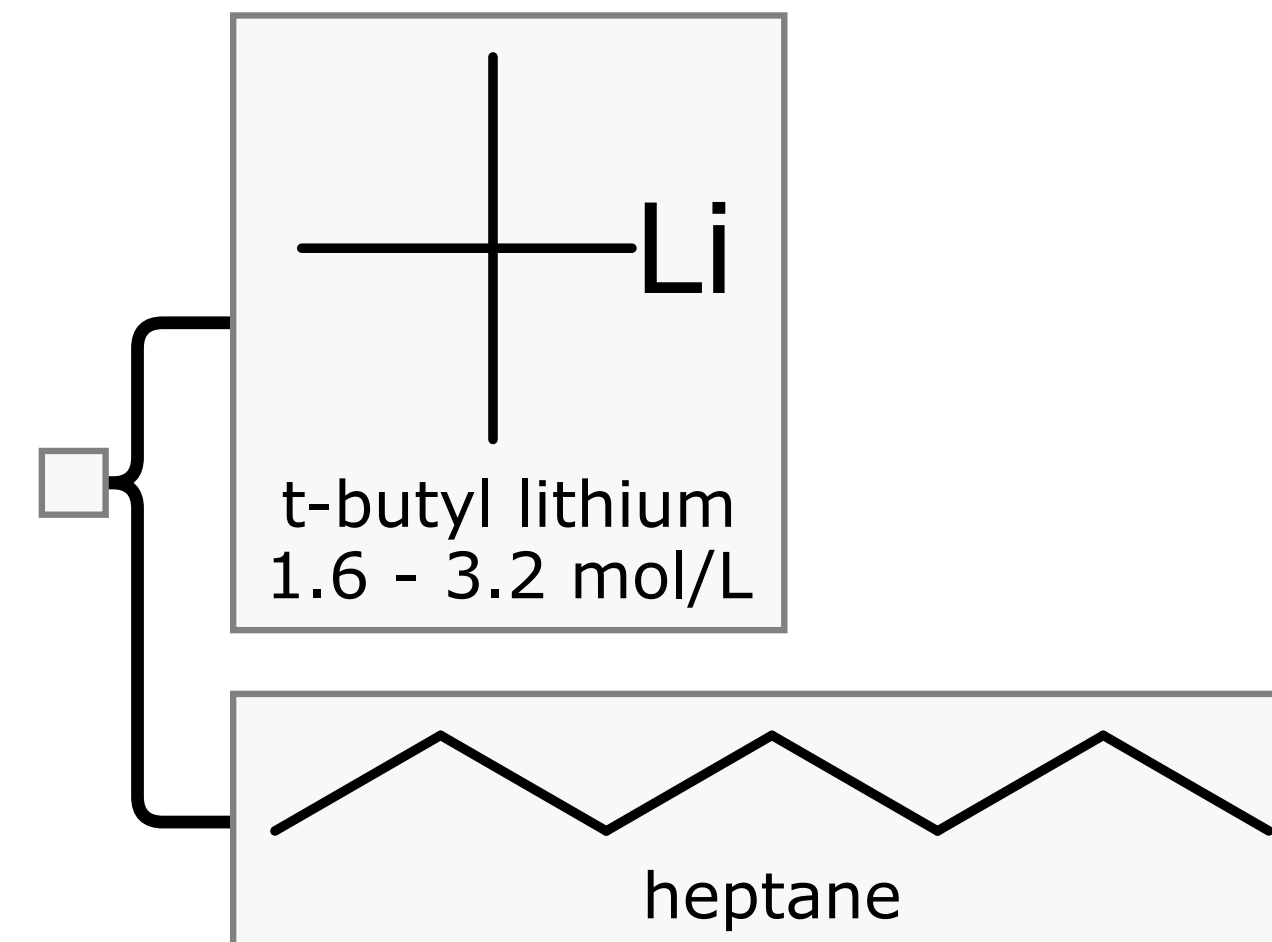
volatile
extreme danger
sealed ampule



storage 2-8 °C
flash point -49 °C



less
dangerous



storage 15 °C
flash point -4 °C

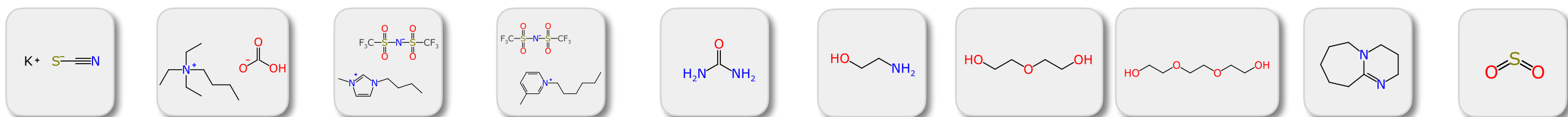
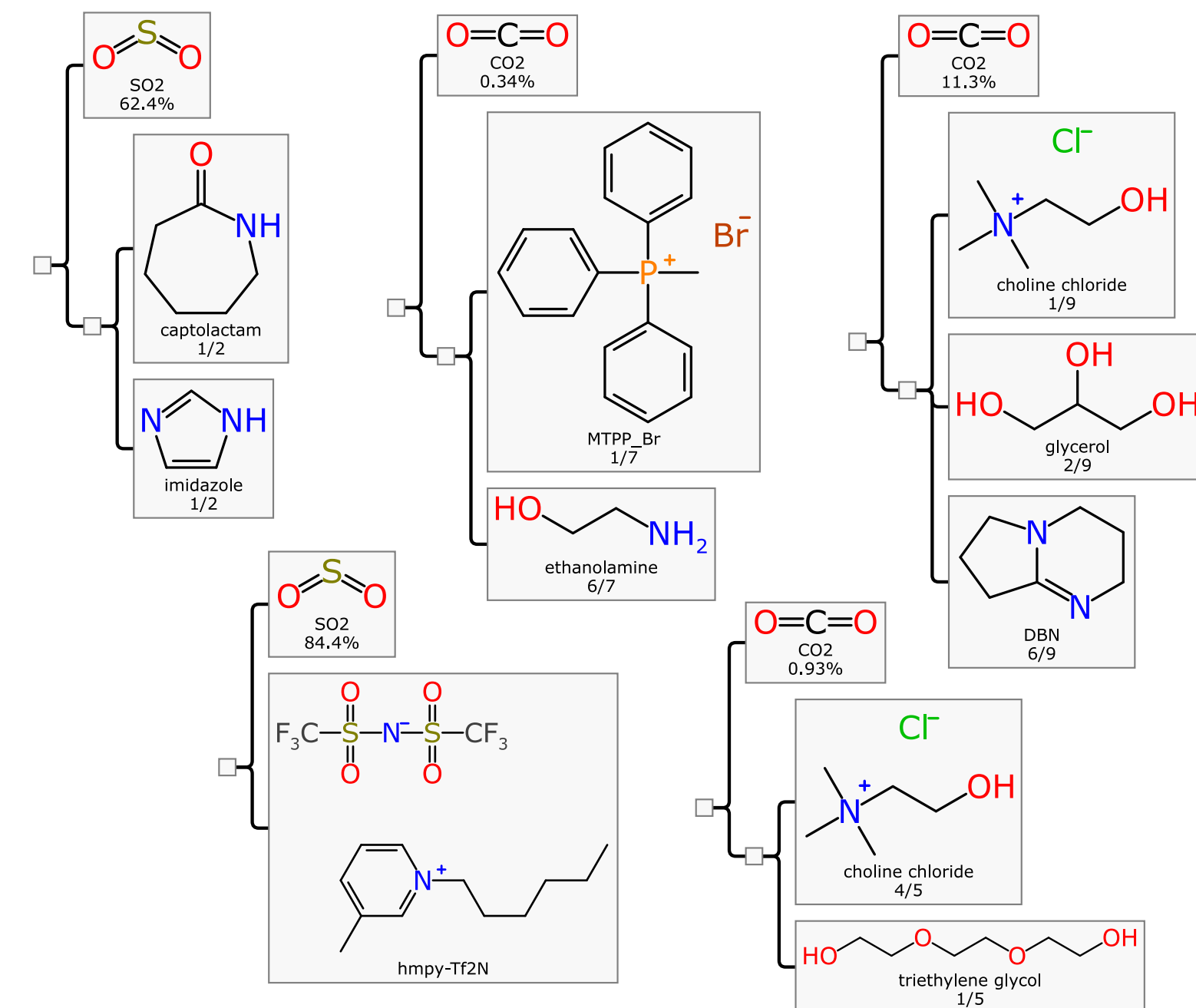
Mixtures QSAR

- ❖ Using chemical structure to model properties has a long history
- ❖ Modelling mixtures adds additional layers:
 - data capture has to be solved first
 - adapt cheminformatics for *multiple* structures
 - factor in concentrations
 - use other descriptors when structures absent
- ❖ Several demo examples...

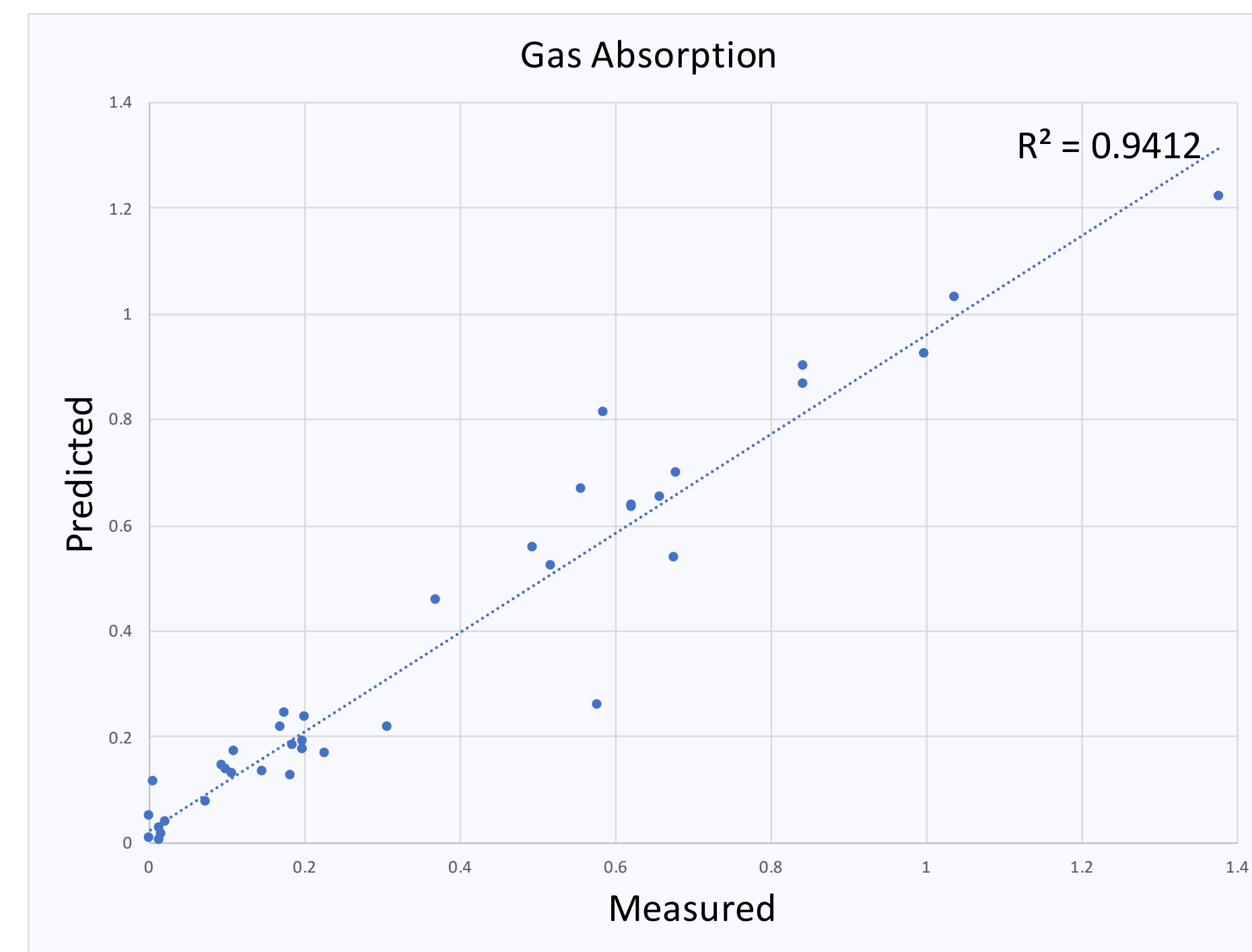
Example 2: Eutectic solubility

✿ <https://bit.ly/3IWFKJV>

✿ Have solubility of CO₂ and SO₂ in various solvent combinations



0	0.01852	0.01389	0.01258	0	0.04762	0.04167	0.03922	0.5455	0
0	0.08	0.01905	0.01667	0.6667	0.07843	0.07018	0.06667	0	0
0	0.05556	0.01961	0.01709	0	0.2963	0.2424	0.2222	0	0
0	0.06122	0.02317	0.02041	0.06593	0.8571	0.2449	0.2286	0	0
0	0.06122	0.02317	0.02041	0.06593	0.8571	0.2449	0.2286	0.004082	0
0	0.06548	0.02317	0.02041	0.06593	0.8571	0.2449	0.2286	0	0
0	0.1429	0.6471	0.1509	0	0.0303	0.02857	0.02778	0.02174	0
0	0.07317	0.4054	0.09615	0	0.03448	0.03226	0.03125	0.02381	0
0	0.1333	0.03902	0.03478	0	0.1882	0.1818	0.2	0	0
0	0.1333	0.03902	0.03478	0	0.1882	0.1684	0.16	0.006452	0
0	0.1333	0.03902	0.03478	0	0.1882	0.2	0.1818	0	0
0	0.08333	0.02439	0.02174	0	0.1176	0.1053	0.1	0	1
0	0.25	0.08163	0.07407	0.1154	0.1481	0.1379	0.1333	0	1
0.25	0.01974	0.0163	0.03	0.03409	0	0	0	0.1406	1
0.25	0.1154	0.04167	0.03659	0.4167	0.05	0	0	0	1
0.1875	0.1154	0.04167	0.03659	0.4167	0.05	0	0	0	1
0.1875	0.01974	0.0163	0.03	0.03409	0	0	0	0.1406	1
0.3	0.072	0.01714	0.015	0.6	0.04615	0	0	0	1
0	0.03226	0.0125	0.01111	0	0.1333	0.4545	0.5	0.5	0
0	0.07778	0.5	0.1633	0.01429	0.01351	0.01282	0.0125	0.5	0
0	1	0.1556	0.14	0.12	0.07143	0.06667	0.06452	0	0
0	0.05556	0.01754	0.0155	0	0.1429	0.125	0.1176	0	0
0	0.2292	0.08108	0.07143	0.02273	0.025	0.02273	0.02174	0.09375	1
0	0.07639	0.02703	0.02381	0.007576	0.008333	0.007576	0.007246	0.03125	1
0	0.07692	0.02778	0.02439	0.2778	0.03333	0	0	0.09375	1
0	0.01316	0.01316	0.02	0.02273	0	0	0	0.09375	1
0	0.1489	0.8788	0.3673	0.02703	0.02564	0.02439	0.02381	0.01923	1
0	0.14	0.3265	1	0.025	0.02381	0.02273	0.02222	0.01818	1
0	0.05926	0.02222	0.01951	0	0.05714	0.05	0.04706	0	0
0	0.08654	0.04286	0.0375	0	0.2727	0.2308	0.2143	0	0
0	0.07759	0.03947	0.03488	0	0.15	0.1324	0.125	0	0
0	0.1429	0.03333	0.02963	0.1333	0.07843	0.07018	0.06667	0	0
0	0.072	0.01951	0.01739	0.6	0.09412	0.08421	0.08	0	0
0	0.08	0.01905	0.01667	0.6667	0.07843	0.07018	0.06667	0	0
0	0.08571	0.02041	0.01786	0.7143	0.06723	0.06015	0.05714	0	0
0	0.08333	0.02439	0.02174	0.25	0.1176	0.1053	0.1	0	0
0	0.08333	0.02439	0.02174	0.3	0.1176	0.1053	0.1	0	0
0	0.08333	0.02439	0.02174	0.35	0.1176	0.1053	0.1	0	0
0	0.08333	0.02439	0.02174	0.25	0.075	0.1053	0.1	0	0
0	0.08333	0.02439	0.02174	0.3	0.05	0.1053	0.1	0	0
0	0.08333	0.02439	0.02174	0.35	0.025	0.1053	0.1	0	0



Polymers: Chi

- ❖ χ is a measure of entropy of mixing
- ❖ Values for polymer + [polymer or solvent or other]
- ❖ 2-component mixtures
- ❖ Data taken from Materials Genome Project...

Materials Genome Project Database Applications

Flory-Huggins Chi (χ) Database

Show 10 entries

Info	Compound 1	Compound 2	χ	χN	Measured at T (K)	Reference
i	poly(lactic acid)	poly(6-methyl- ϵ -caprolactone)	$-0.1 + \frac{61.2}{T}$		428.15	10.1021/ma201063t 10.1021/ma201063t
i	poly(methyl acrylate)	silicon dioxide	-0.3		423.15	10.1021/ma200205j
i	poly(methyl acrylate)	2-butanone	-0.45		423.15	10.1021/ma200205j 10.1021/ef502918t 10.1021/ma035830m
i	water	poly(ethylene oxide)	$0.8 - \frac{100}{T}$		298.15	10.1021/ma301193h
i	polythiophene	tetrahydrofuran	3.34		300	10.1021/ma401345n
i	poly(3-methylthiophene)	tetrahydrofuran	11.43		300	10.1021/ma401345n
i	poly(3-butylthiophene)	tetrahydrofuran	7.41		300	10.1021/ma401345n
i	poly(3-hexylthiophene)	tetrahydrofuran	0.39		300	10.1021/ma401345n
i	poly(3-octylthiophene)	tetrahydrofuran	-1.06		300	10.1021/ma401345n
i	poly(3-dodecylthiophene)	tetrahydrofuran	-0.38		300	10.1021/ma401345n

Showing 1 to 10 of 263 entries Previous 1 2 3 4 5 ... 27 Next

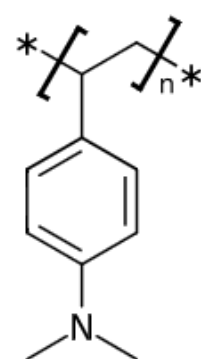
Spreadsheet

- ❖ Content arrives in tabular form
- ❖ Expanding into full mixtures - 2 steps:
 1. register each named structure
 2. use Vault to compose into mixtures

The screenshot shows a Microsoft Excel spreadsheet with the following data:

	A	B	C	D
1	ID	Component1	Component2	Chi
2	PolyMix001	poly(lactic acid)	poly(6-methyl-ε-caprolact	0.0429
3	PolyMix002	poly(methyl acrylate)	silicon dioxide	-0.3000
4	PolyMix003	poly(methyl acrylate)	2-butanone	-0.4500
5	PolyMix004	water	poly(ethylene oxide)	0.4646
6	PolyMix005	polythiophene	tetrahydrofuran	3.3400
7	PolyMix006	poly(3-methylthiophene)	tetrahydrofuran	11.4300
8	PolyMix007	poly(3-butylthiophene)	tetrahydrofuran	7.4100
9	PolyMix008	poly(3-hexylthiophene)	tetrahydrofuran	0.3900
10	PolyMix009	poly(3-octylthiophene)	tetrahydrofuran	-1.0600
11	PolyMix010	poly(3-dodecylthiophene)	tetrahydrofuran	-0.3800
12	PolyMix011	polystyrene	poly(4-hydroxystyrene)	0.6800
13	PolyMix012	poly(ethylene oxide)	poly(methyl acrylate)	-0.0400
14	PolyMix013	polybutadiene	tetrahydrofuran	0.3990
15	PolyMix014	polybutadiene	tetrahydrofuran	0.4030
16	PolyMix015	polybutadiene	tetrahydrofuran	0.4000
17	PolyMix016	polybutadiene	tetrahydrofuran	0.4000
18	PolyMix017	polybutadiene	tetrahydrofuran	0.4110
19	PolyMix018	polybutadiene	tetrahydrofuran	0.4200
20	PolyMix019	polyisoprene	cyclohexane	0.3940
21	PolyMix020	polyisoprene	cyclohexane	0.3730
22	PolyMix021	polyisoprene	cyclohexane	0.3950
23	PolyMix022	polyisoprene	cyclohexane	0.4370
24	PolyMix023	polyisoprene	cyclohexane	0.4380
25	PolyMix024	polyisoprene	cyclohexane	0.4400
26	PolyMix025	poly(4-tert-butylstyrene)	poly(methyl methacrylate	0.0966

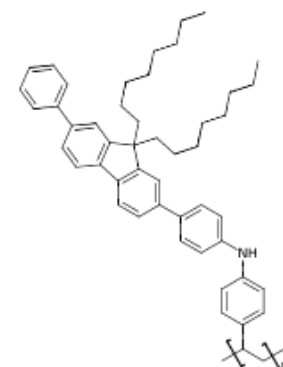
Register structures



poly[N-(4-vinylbenzyl)-
N,N-diethylamine]

Synonyms: *(no synonyms)*

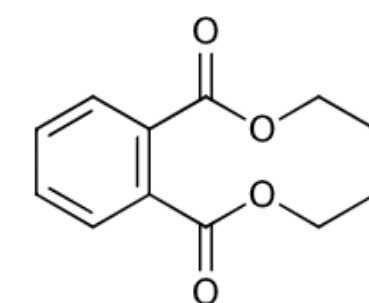
Protocols: 0



poly[(9,9-dioctylfluorenyl-
2,7-diyl)-co-5,5-(4',7'-di-2-
thienyl-2',1',3'-
benzothiadiazole)]

Synonyms: *(no synonyms)*

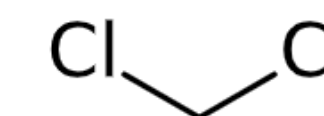
Protocols: 0



diethyl phthalate

Synonyms: *(no synonyms)*

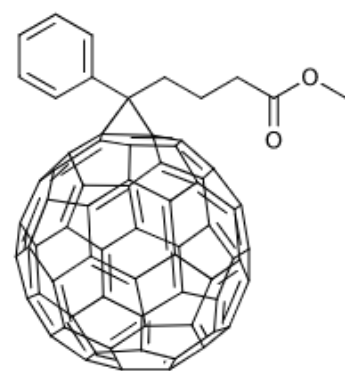
Protocols: 0



dichloromethane

Synonyms: *(no synonyms)*

Protocols: 0

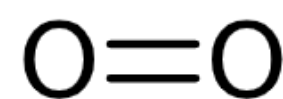


phenyl-C61-butyric acid
methyl ester

Synonyms: [6,6]-phenyl-C61

butyric acid methyl ester

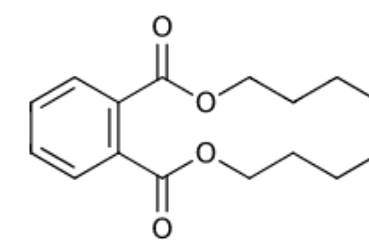
Protocols: 0



oxygen

Synonyms: *(no synonyms)*

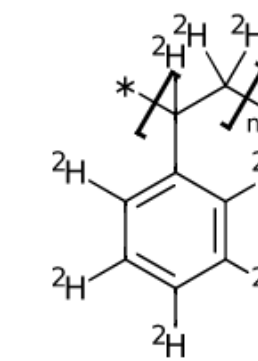
Protocols: 0



dibutyl phthalate

Synonyms: *(no synonyms)*

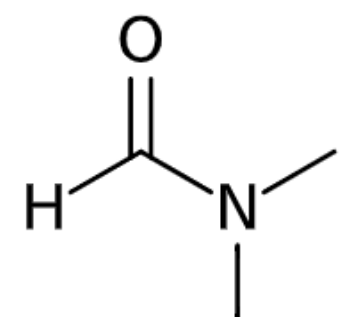
Protocols: 0



deuterated polystyrene

Synonyms: *(no synonyms)*

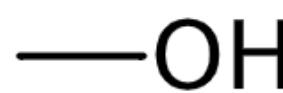
Protocols: 0



n,n-dimethylformide

Synonyms: *(no synonyms)*

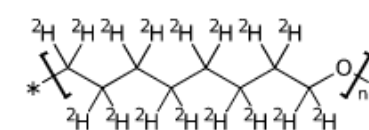
Protocols: 0



methanol

Synonyms: *(no synonyms)*

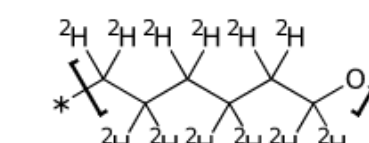
Protocols: 0



deuterated poly(octylene
oxide)

Synonyms: *(no synonyms)*

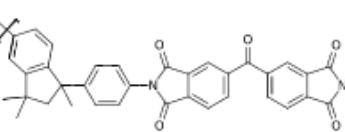
Protocols: 0



deuterated poly(hexene
oxide)

Synonyms: *(no synonyms)*

Protocols: 0



matrimid

Synonyms: *(no synonyms)*

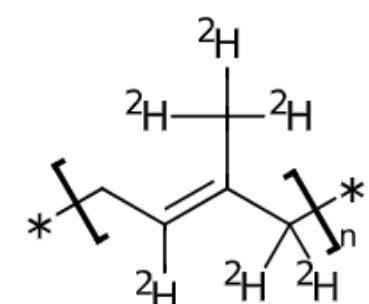
Protocols: 0



ethylene dichloride

Synonyms: 1,2-dichloroethane

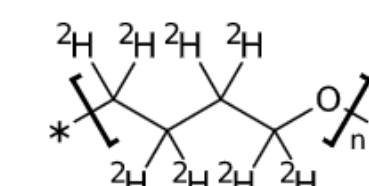
Protocols: 0



deuterated poly(ethylene-
alt-propylene)

Synonyms: *(no synonyms)*

Protocols: 0



deuterated poly(butylene
oxide)

Synonyms: *(no synonyms)*

Protocols: 0

Markup into mixtures

❖ CDD Vault importer: *Compose Mixtures*

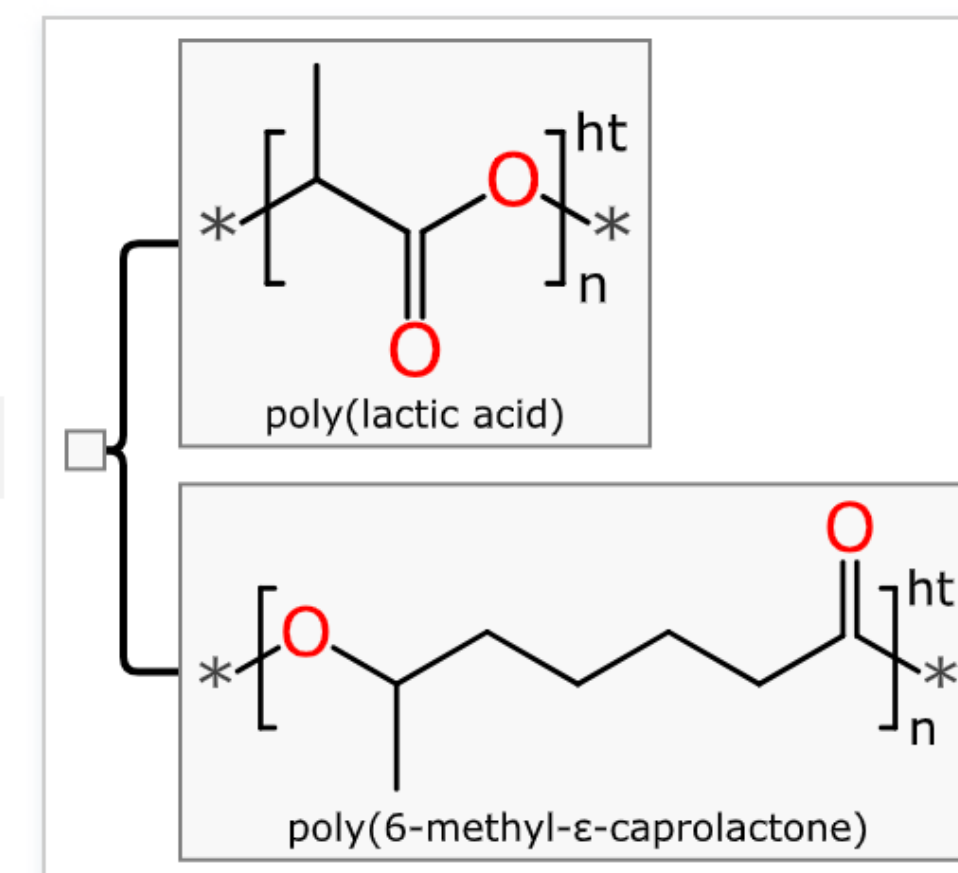
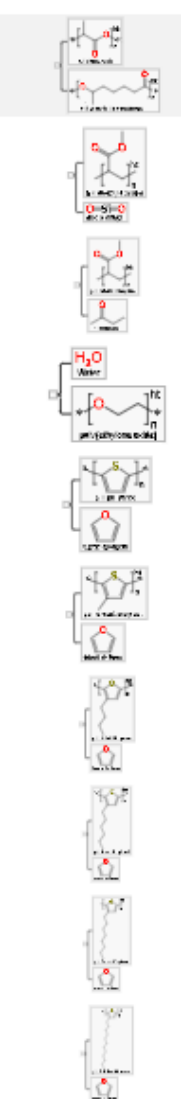
- Not part of mixture
- Chemical structure [Molfile]
- Chemical structure [SMILES]
- Molecule lookup by name
- CDD Vault identifier
- Component name (text)
- Quantity
- Units

Component

1	2	3
4	5	6
7	8	9
	<input type="text"/>	?

	A	B ⇒ 1 Vault identifier	C ⇒ 2 Vault identifier	D	E	
	ID	Component1	Component2	Chi	Chi-Raw	Chi-N
1	PolyMix001	poly(lactic acid)	poly(6-methyl-ε-caprolactone)	0.04294	-0.1 + 61.2T	
2	PolyMix002	poly(methyl acrylate)	silicon dioxide	-0.30000	-0.3	
3	PolyMix003	poly(methyl acrylate)	2-butanone	-0.45000	-0.45	
4	PolyMix004	water	poly(ethylene oxide)	0.46460	0.8 - 100T	
5	PolyMix005	polythiophene	tetrahydrofuran	3.34000	3.34	
6	PolyMix006	poly(3-methylthiophene)	tetrahydrofuran	11.43000	11.43	
7	PolyMix007	poly(3-butylthiophene)	tetrahydrofuran	7.41000	7.41	
8	PolyMix008	poly(3-hexylthiophene)	tetrahydrofuran	0.39000	0.39	
9	PolyMix009	poly(3-octylthiophene)	tetrahydrofuran	-1.06000	-1.06	
10	PolyMix010	poly(3-dodecylthiophene)	tetrahydrofuran	-0.38000	-0.38	

and 237 additional rows



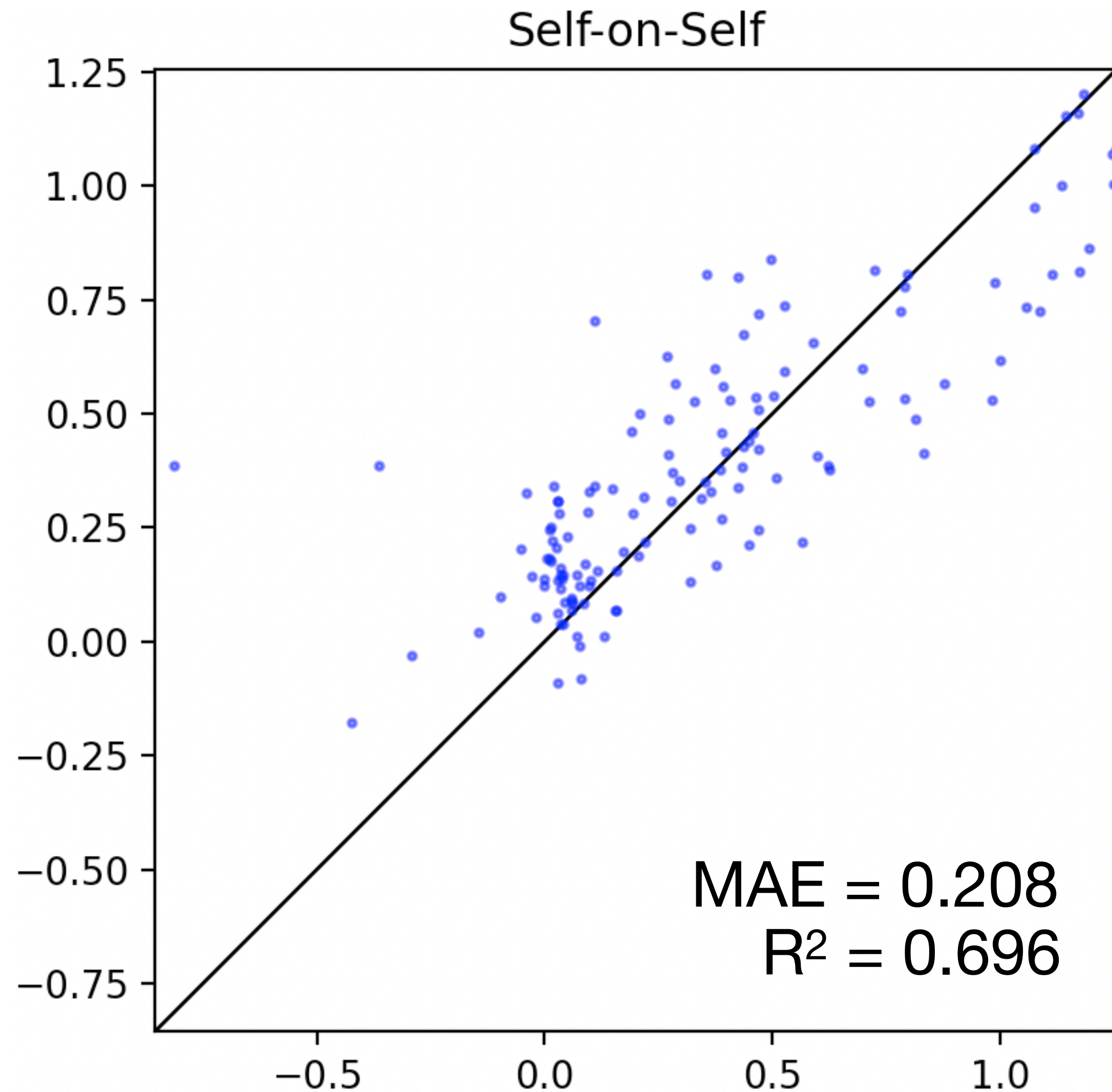
Gather data

- ❖ Select all chemical objects with a *Flory-Huggins Chi* property
- ❖ 172 unique datapoints
- ❖ Data gathering is generic: applies to any mixture-formatted content with this property

172 Selected: [Plot](#) [Export](#) [Add to collection](#) [Build model](#) [Flag outliers](#) [Customize your report](#) [Save this search](#)

Select...	Molecule	Flory-Huggins
<input checked="" type="checkbox"/>	<p>polyisoprene poly(lactic acid)</p> <p>PolyMix245 McKerrow Vault</p>	0.799 0.689 0.419
<input checked="" type="checkbox"/>	<p>poly(propylene oxide) thymine</p> <p>PolyMix244 McKerrow Vault</p>	10.6
<input checked="" type="checkbox"/>	<p>polyethylene 2,6-diaminotriazine</p> <p>PolyMix243 McKerrow Vault</p>	3.00
		8.50

Results



- ❖ Can propose:
 - ▶ arbitrary polymer structure
 - ▶ arbitrary other molecule
 - ▶ predict entropy of mixing (χ)
- ❖ Dataset is small and diverse...
- ❖ ... cross validation not great, need more data
- ❖ Database of mixtures & χ -values?
- ❖ Bring into training set

The Future

- ❖ Digitise mixtures, like for molecules - *self describing*
- ❖ Put them in repositories
 - with measured properties
 - open databases (like *PubChem*)
 - or just a huge data lake
- ❖ Then focus on
 - queries, analysis, model building - the fun part!
 - describing components that are not simple molecules

Acknowledgements

- ♣ Leah McEwen
- ♣ InChI Trust & IUPAC
- ♣ NIH Grant 1R43TR002528-01
- ♣ The CDD Vault team



alex@collaborativedrug.com

(PS: we're hiring)

