

Canonical Line Notations

InChI vs SMILES

Krisztina Boda



- Compound naming
- InChI
- SMILES
- Molecular equivalency
 - Isomorphism
 - Kekule
 - Tautomers
- Finding duplicates

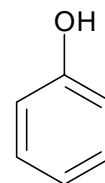


What's Your Name?

1. Unique numbers

- CAS registry number
- BRN Beilstein Registry Number
- EC (European commission) number
- CID (PubChem compound id)

one-to-one

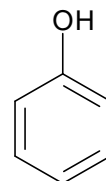


— CAS: 108-95-2

2. Chemical Nomenclatures

- IUPAC name(s)
- Traditional name(s)
[carbolic acid, benzenol, phenylic acid, std.]

one-to-many



IUPAC [ENG] phenol

IUPAC [HUN] fenol

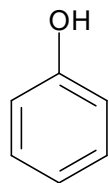


What's Your Name?

3. Hashed identifiers

- NCI/CADD structure identifiers
- InChIKey

many-to-one

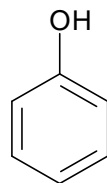


InChIKey= ISWSIDIOOBBQZ-UHFFFAOYSA-N

4. Chemical line notations

- WLN (Wiswesser Line Notation)
- **SMILES**
- **InChI**

one-to-one ?



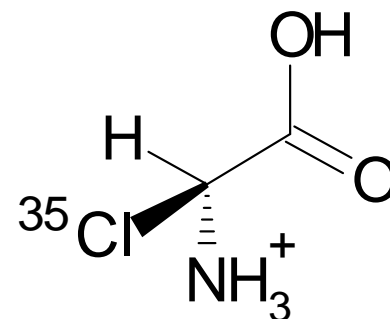
InChI=1S/C6H6O/c7-6-4-1-3-5-6/h1-5,7H

c1ccccc1O



InChI Identifier (**I**nternational **C**hemical **I**dentifier)

- Developed by IUPAC and NIST [first release 2005]
- One reference implementation
- Multiple hierarchical layers
 - No bond order
- Not a file format!



main layer



InChI=1S/C2H4ClNO2/c3-1(4)2(5)6/h1H,4H2,(H,5,6)/p+1/t1-/m1/s1/i3+0

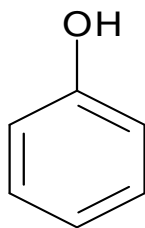
version

charge

stereo

isotope

SMILES (Simplified Molecular Input Line Entry Specification)



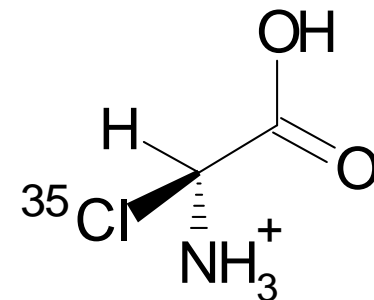
```

c1ccccc1O
Oc1ccccc1
C1=CC=CC=C1O
C1=CC=C(C=C1)O
[CH]1=[CH][CH]=[CH][CH]=C1[OH]
  
```

- Canonical SMILES [OEChem: c1ccc(cc1)O]
 - Unique name for each molecule in one system
 - Not a global identifier
- Canonical Isomeric SMILES
 - Encode isotope, double bond and chiral configuration

```

C(C(=O)O)([NH3+])Cl
[C@H](C(=O)O)([NH3+])[35Cl]
  
```

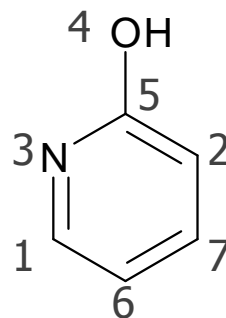


[1] SMILES, a chemical language and information system. 1. Introduction to methodology and encoding rules
David Weininger, J. Chem. Inf. Comput. Sci., **1988** (28)

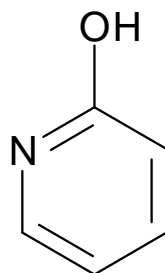
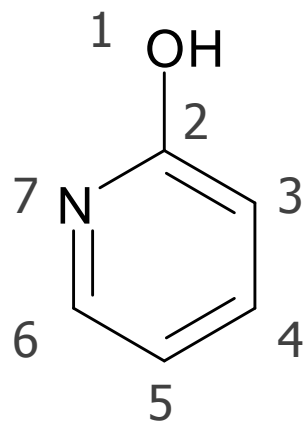


Molecule Registration

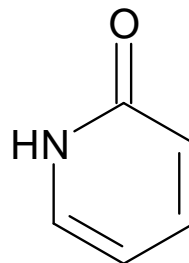
Is this already
in my
database?



Isomorph



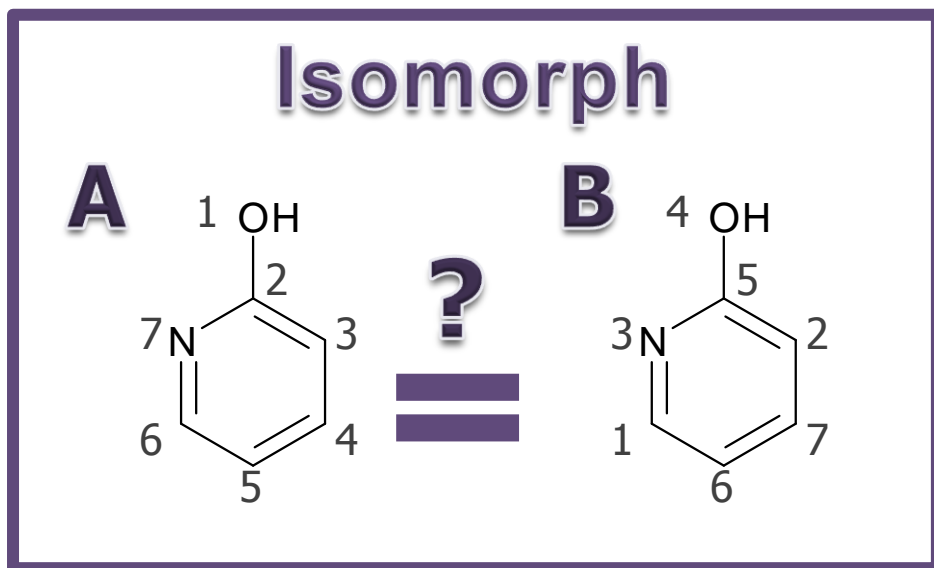
Kekule



Tautomer



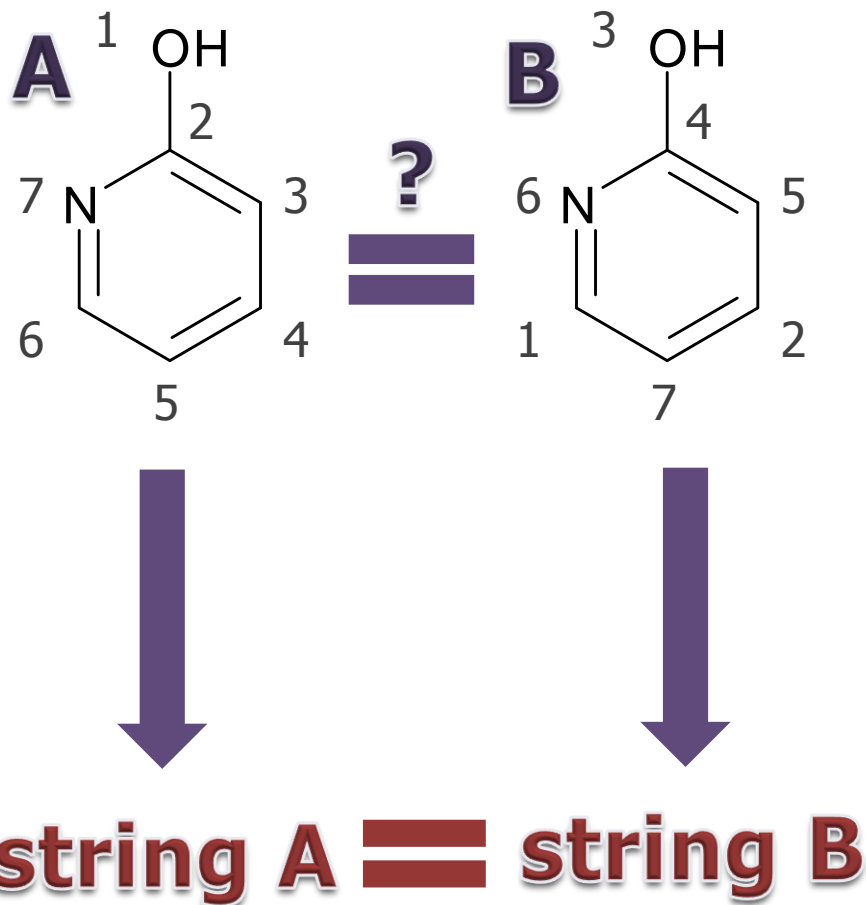
Molecule Equivalency - Isomorphism



- NP problem
- Complexity = $HvyAtm!$
- Impractical for molecule registration!



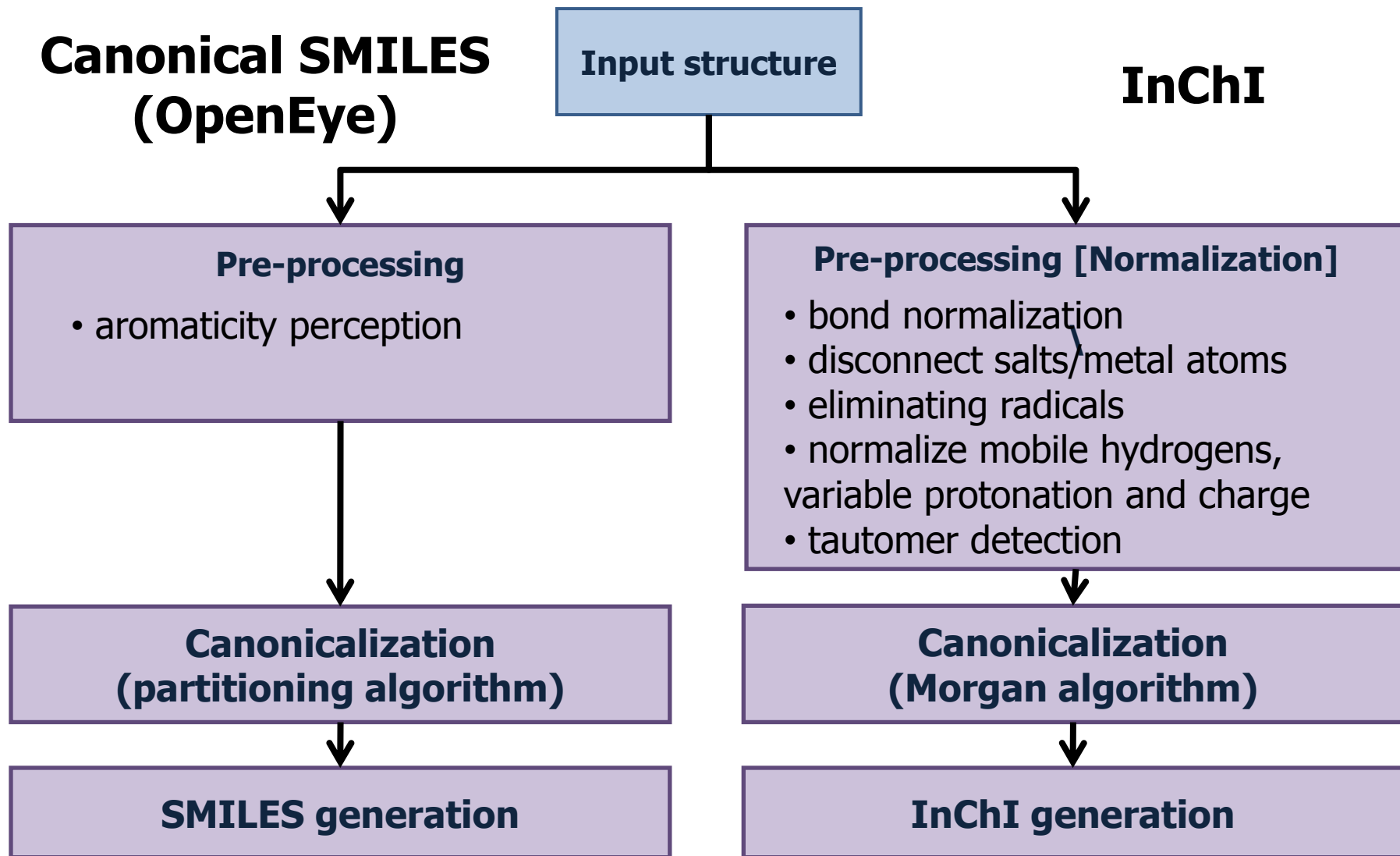
Canonicalization

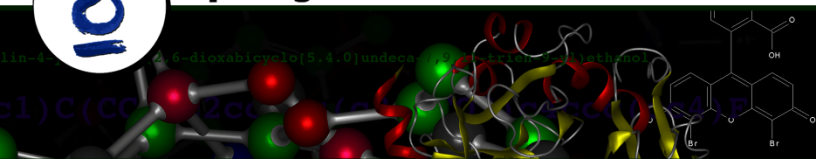


- Transforms any connection table into a unique canonical form (independent of atom indices)
- Comparing molecules \Rightarrow comparing unique, unambiguous string representations



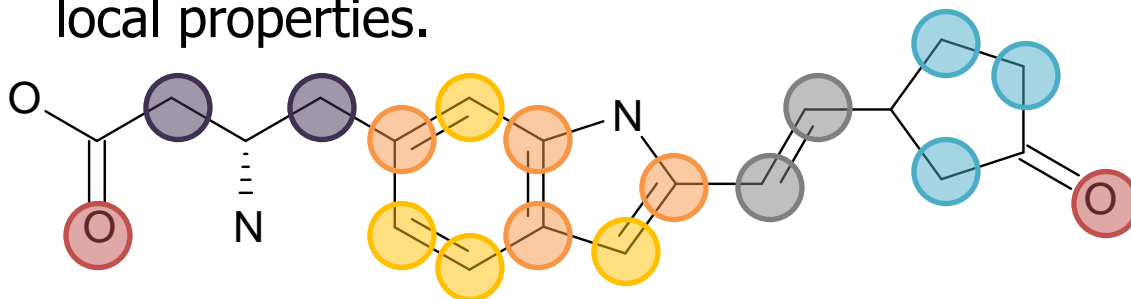
SMILES/InChI Generation





Finding unique atom ordering

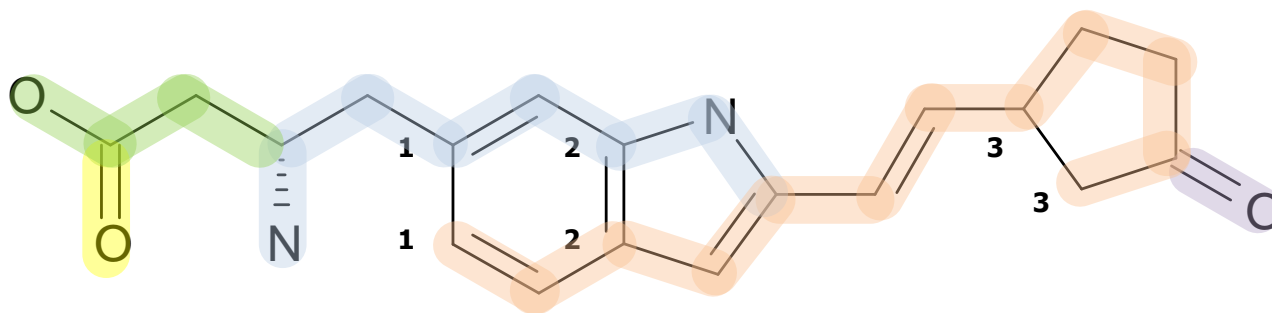
1. Initial partitioning of atoms by distinguishing them by their local properties.



2. Recursively refining atom classes by distinguishing atoms by their neighbors.
3. When each atom belong to a unique partition \Rightarrow unique canonical atom order exists.



SMILES - Generation



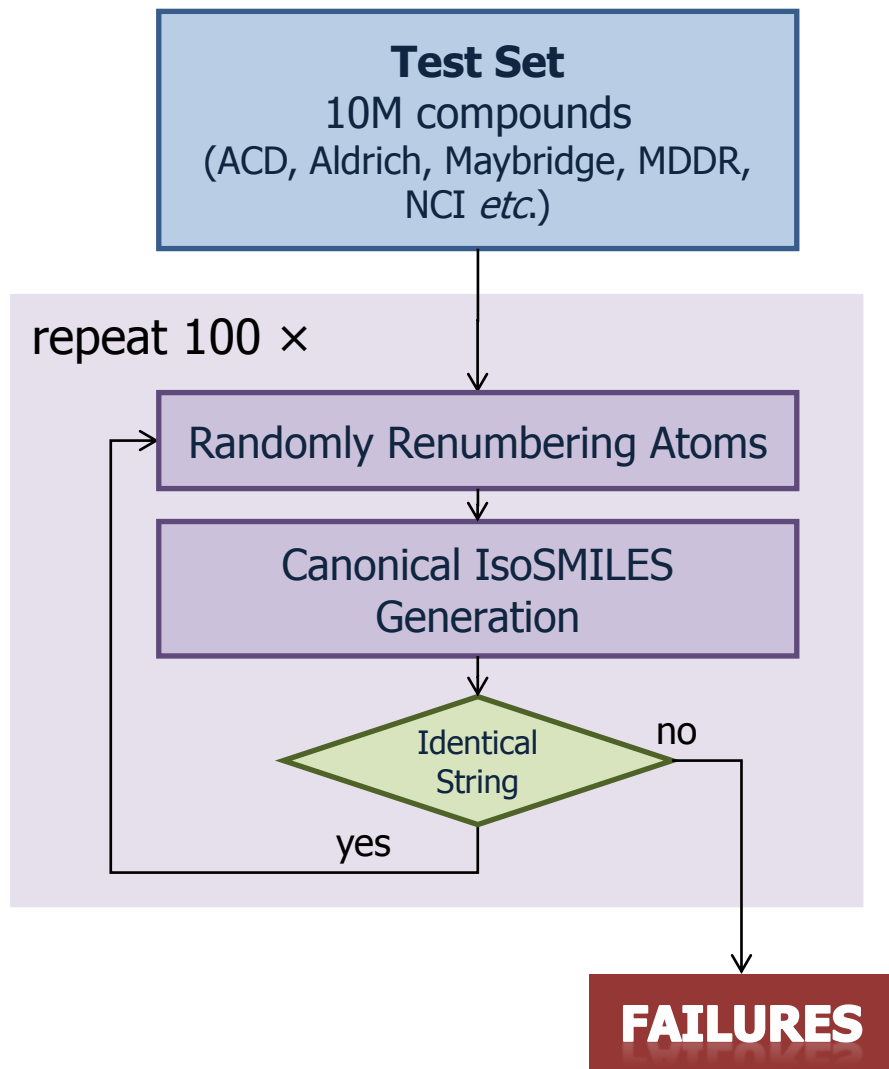
```
c1cc2cc([nH]c2cc1C[C@@H](CC(=O)O)N)/C=C/C3CCCC(=O)C3
```



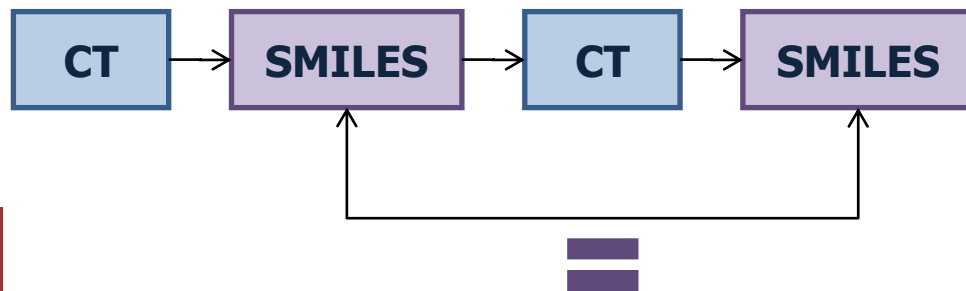
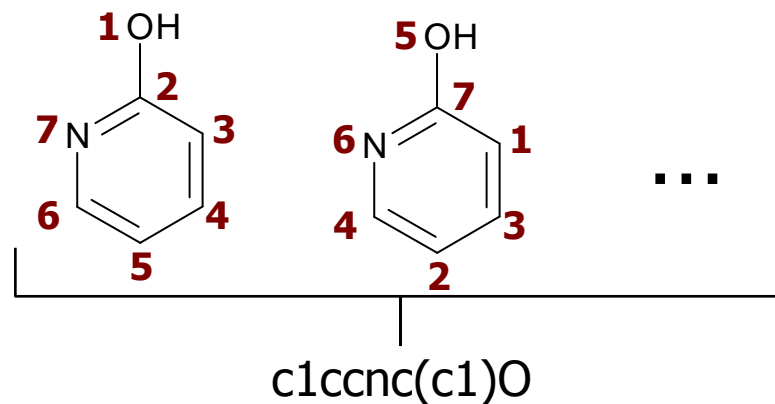
Based on the canonical ordering of atoms, a unique SMILES is generated by depth-first search (DFS).



SMILES Validation



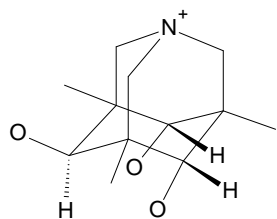
Name has to be independent of atom numbering!





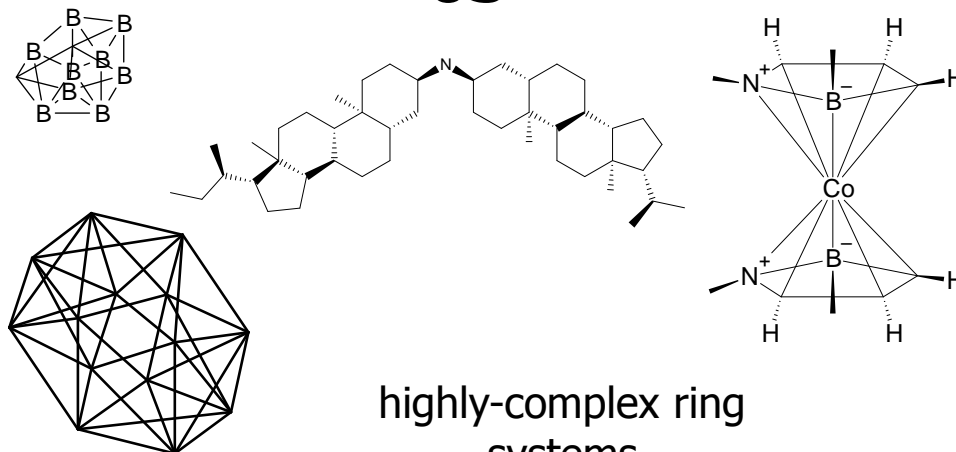
SMILES Validation

7



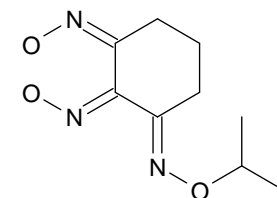
adamantane-like ring systems

63



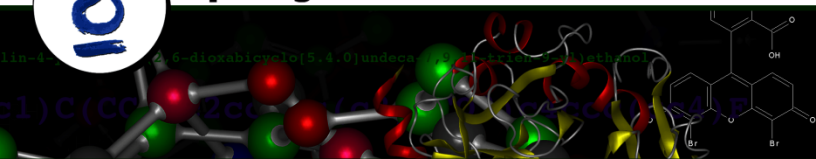
highly-complex ring systems

8



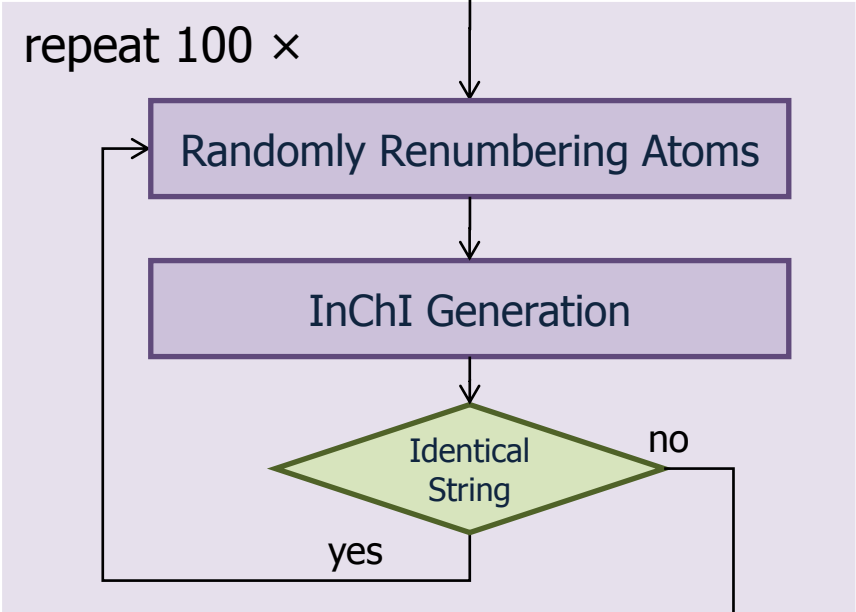
SMILES bond stereo representation problem

Number of total failures: 78 (0.0008%)

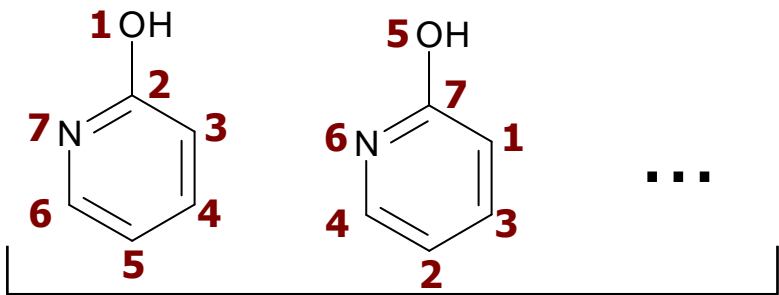


InChI Validation

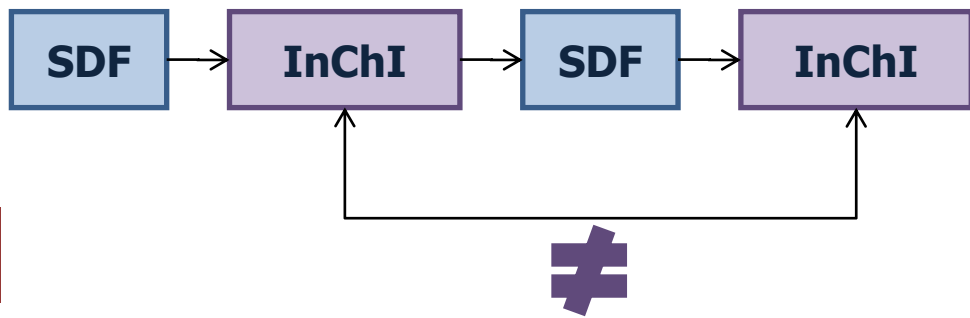
Test Set
10M compounds
(ACD, Aldrich, Maybridge, MDDR,
NCI *etc.*)



Name has to be independent of atom numbering!

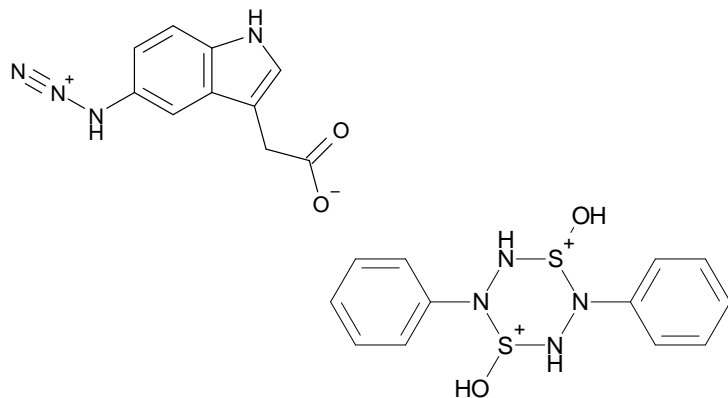


InChI=1S/C5H5NO/c7-5-3-1-2-4-5/h1-4H,(H,6,7)



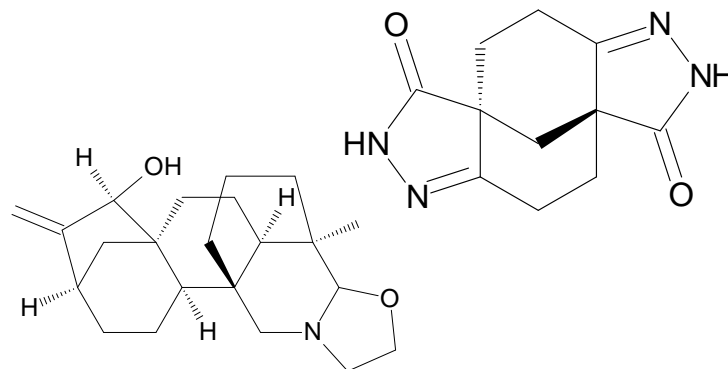
InChI Validation

28



charge layer

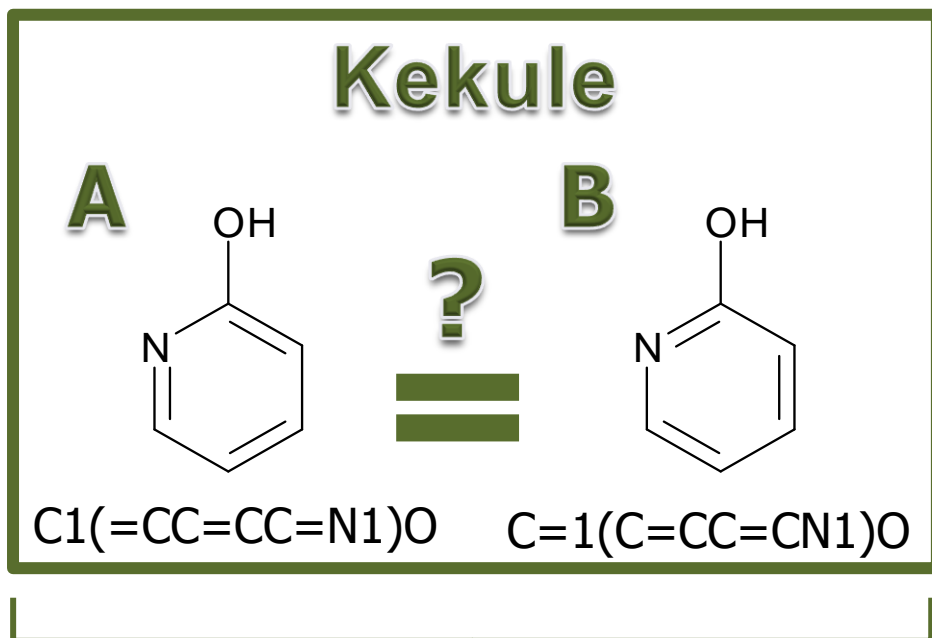
27



stereo layer

Number of total failures: 55 (0.0006%)

Molecule Equivalency (Kekule)



c1cnc(O)ccc1

InChI=1S/c5H5NO/c7-5-3-1-2-4-6-5/h1-4H,(H,6,7)

Canonical SMILES

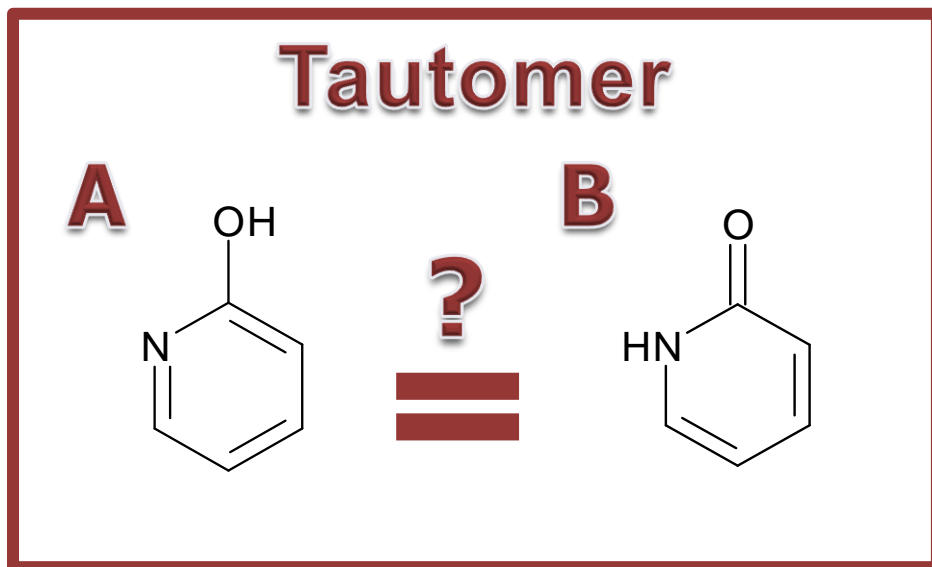
- aromatic form (no bond between aromatic atoms)

InChI

- no bond information



Molecule Equivalency - Tautomer



- J. Comput.-Aided Mol. Des. (June 2010)

For registration:

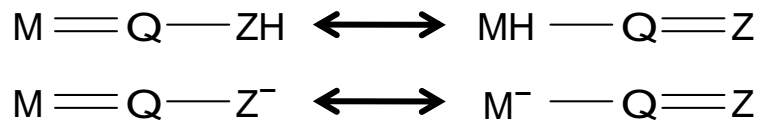
Canonical tautomer does **not have** to be

- the chemically most reasonable
- the lowest in energy

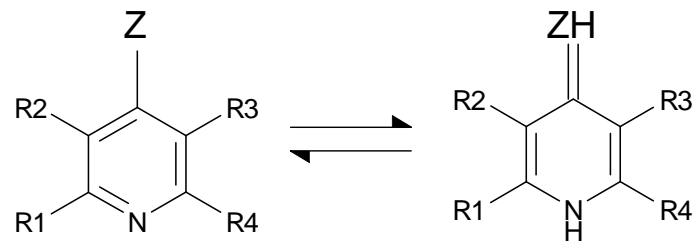
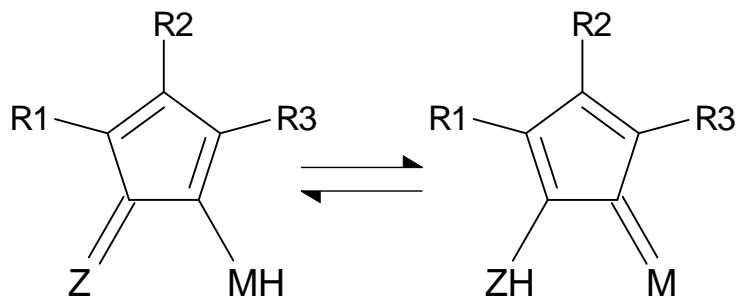
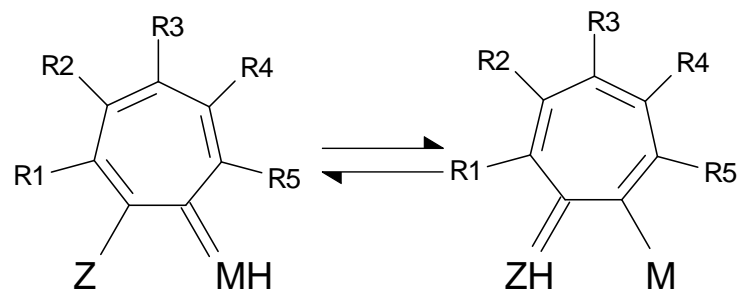
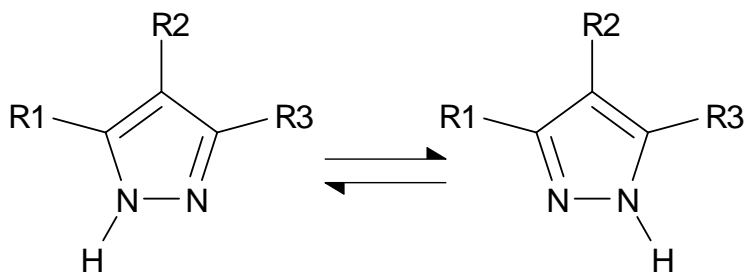
Canonical tautomer **has to be** identical for every tautomer



InChI- Tautomer Detection

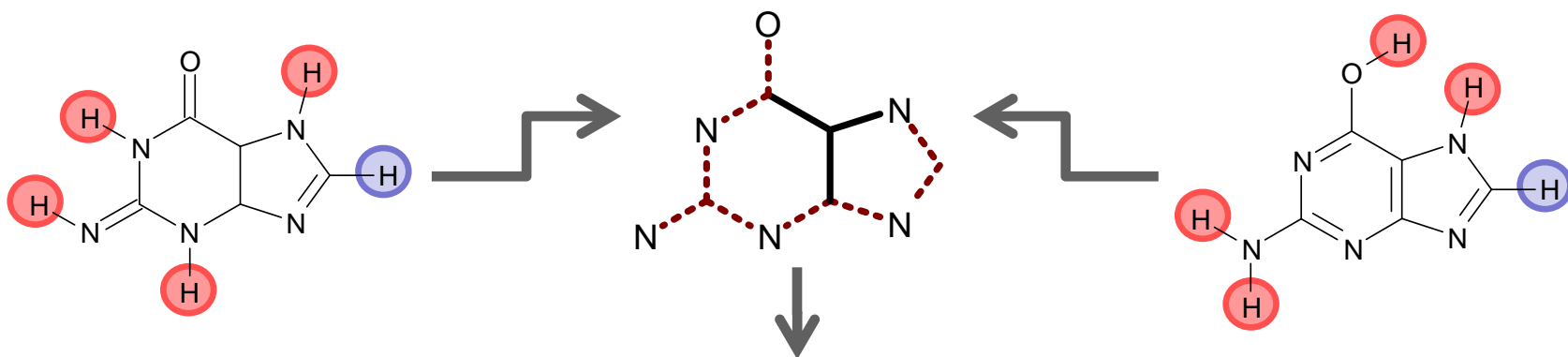


M,Z = trivalent N, bivalent O, S, Se, Te
Q = C, N, S, P, Sb, As, Se, Te, Br, Cl, I
H = H, D, T

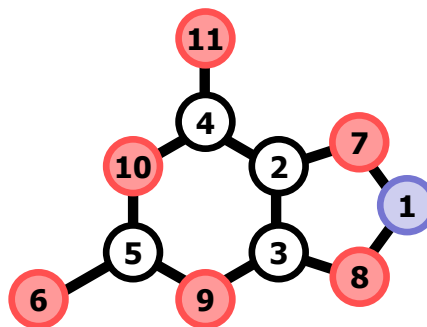




InChI - Handling Tautomers



Same InChI
for all 15
guanine
tautomers!



InChI = 1S/C5H5N5O/c6-5-93-2(4(11)10-5)7-1-8-3/h1H,(H4,6,7,8,9,10,11)



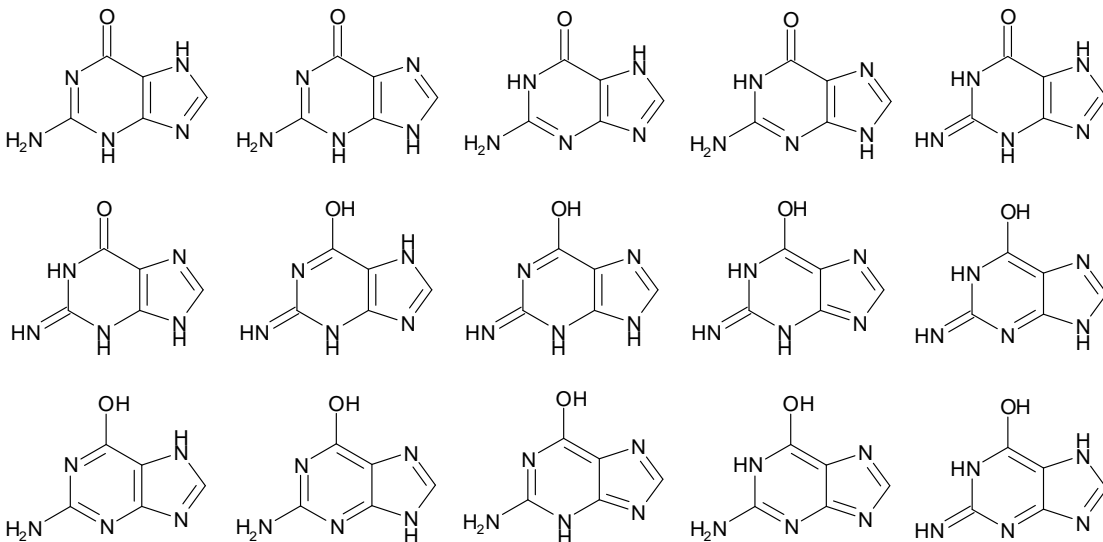
hydrogen layer

SMILES – Handling tautomers

- Identifying and handling tautomers is not part of the language
- Tautomeric structures are explicitly specified
 - No tautomeric bond
 - No mobile hydrogens

```

c1[nH]c2c(n1)[nH]c(nc2=O)N
c1[nH]c2c(=O)[nH]c(nc2n1)N
c1[nH]c2c(n1)c(=O)nc([nH]2)N
c1[nH]c2c(n1)c(=O)[nH]c(n2)N
c1[nH]c2c(n1)[nH]c(=N)[nH]c2=O
c1[nH]c2c(n1)c(=O)[nH]c(=N)[nH]2
c1[nH]c2c(n1)nc(nc2O)N
c1nc-2c(nc([nH]c2n1)N)O
c1[nH]c2c(n1)c(nc(n2)N)O
c1nc-2c([nH]c(nc2n1)N)O
c1[nH]c2c(n1)[nH]c(=N)nc2O
c1[nH]c2c([nH]c(=N)nc2n1)O
c1[nH]c2c(n1)c(nc(=N)[nH]2)O
c1nc-2c([nH]c(=N)[nH]c2n1)O
c1[nH]c2c(n1)c([nH]c(=N)n2)O
  
```





SMILES/InChI Generation

Canonical SMILES (OpenEye)

Input structure

InChI

Pre-processing

- aromaticity perception

QuacPac

Canonicalization (partitioning algorithm)

SMILES generation

Pre-processing [Normalization]

- bond normalization
- disconnect salts/metal atoms
- eliminating radicals
- normalize mobile hydrogens, variable protonation and charge
- tautomer detection

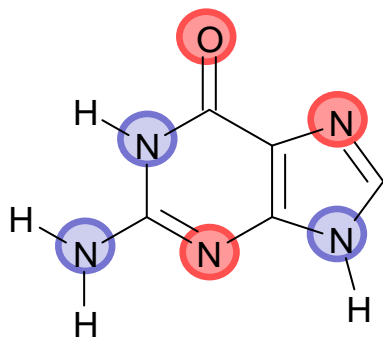
Canonicalization (Morgan algorithm)

InChI generation

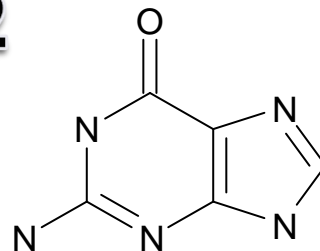


QuacPac – Canonical Tautomer

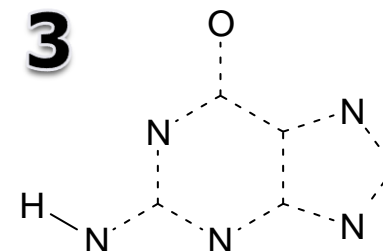
1



2



3



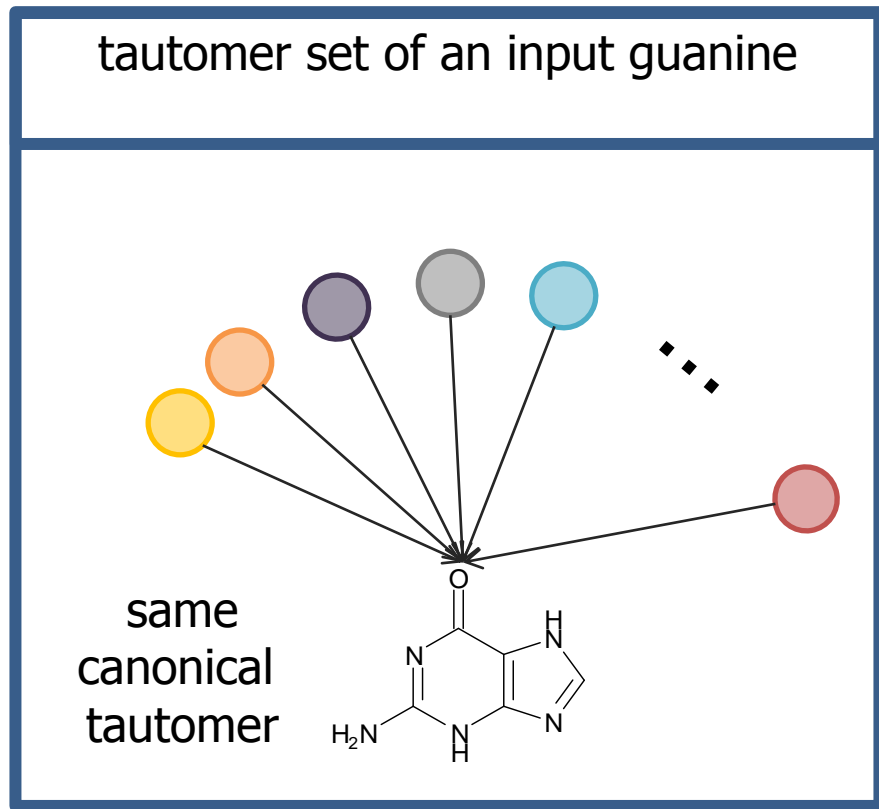
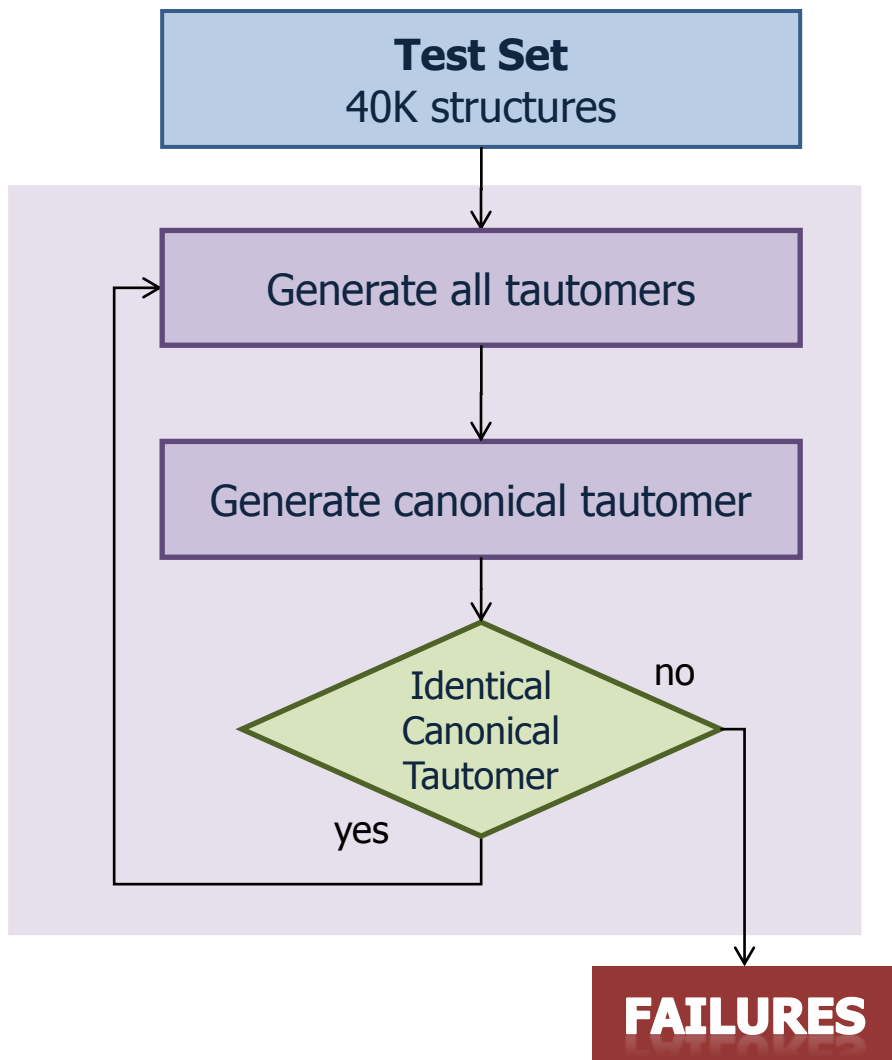
 donor  acceptor

Global algorithm that systematically generates **all** tautomers!

1. atom typing
2. removing (+counting) hydrogens
3. unspecifying bond types
4. enumerating tautomers (DFS)
 - protonating/deprotonating atom
5. Empirical, rule-based scoring to find the canonical tautomer



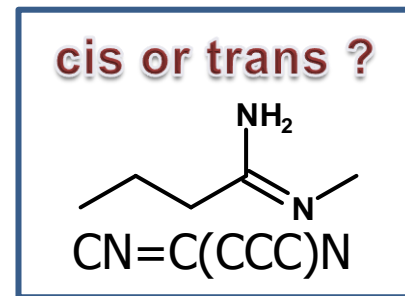
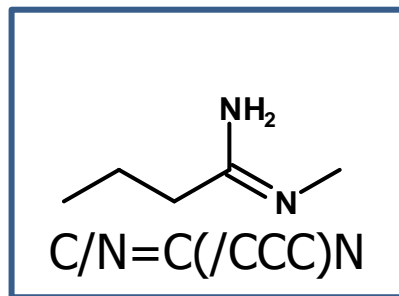
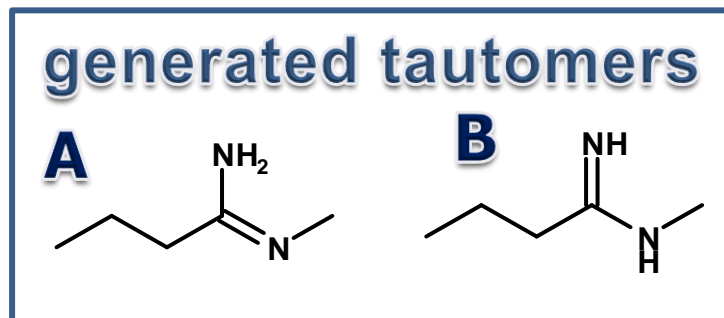
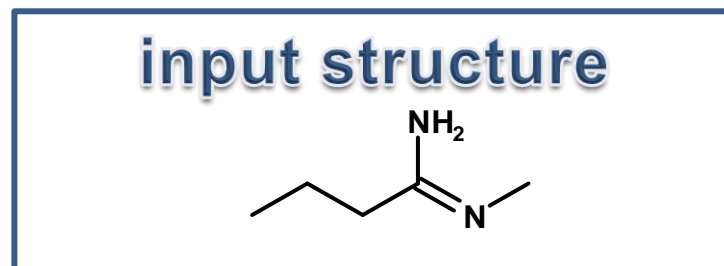
Validation – Canonical Tautomer





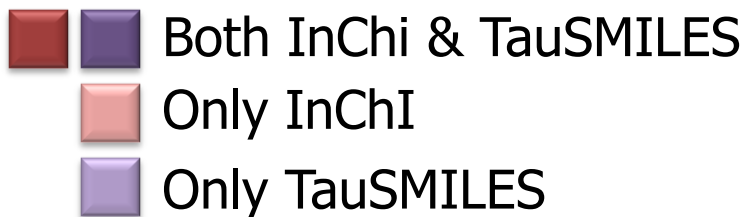
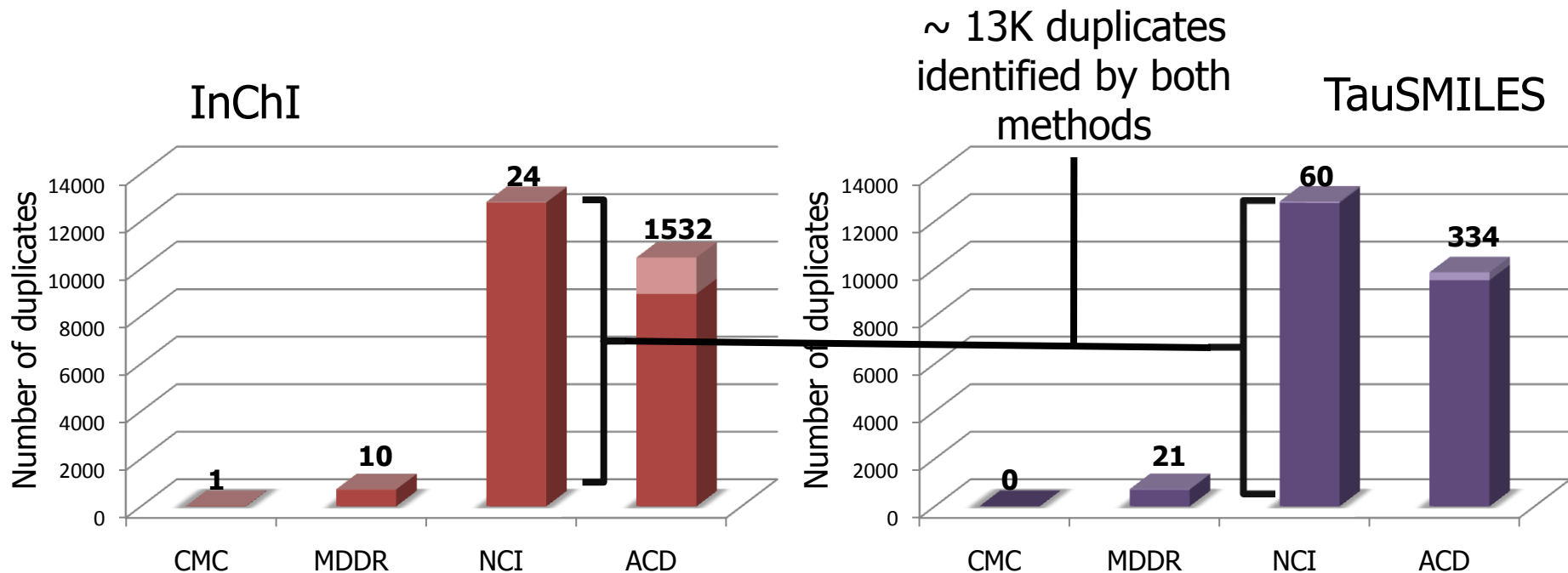
Tautomer Generation for SMILES

- Input: 40K
- Generated tautomers: 230K
- Failures: 272 (0.007%)
- Failures due to changing E/Z geometry 269





Identifying Duplicate Compounds

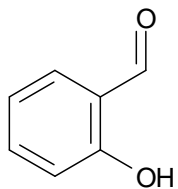


	CMC	MDDR	NCI	ACD
Both InChi & TauSMILES	14	704	12784	8940
Only InChi	1	10	24	1532
Only TauSMILES	0	21	60	334



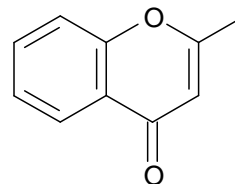
Duplicates

Isomorphs

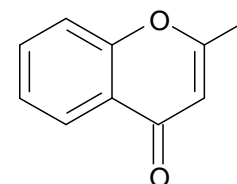


NSC numbers:
49178,51705,83559 ,83560,83561,
83562,97202,112278,187662

Kekules

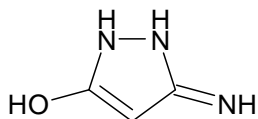


NSC numbers:
74869

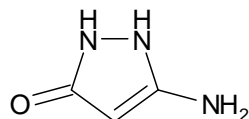


NSC numbers:
660048

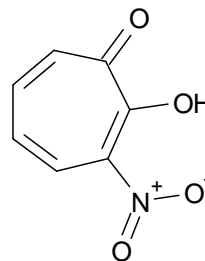
Tautomers



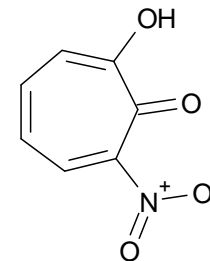
NSC numbers:
60188



NSC numbers:
408479

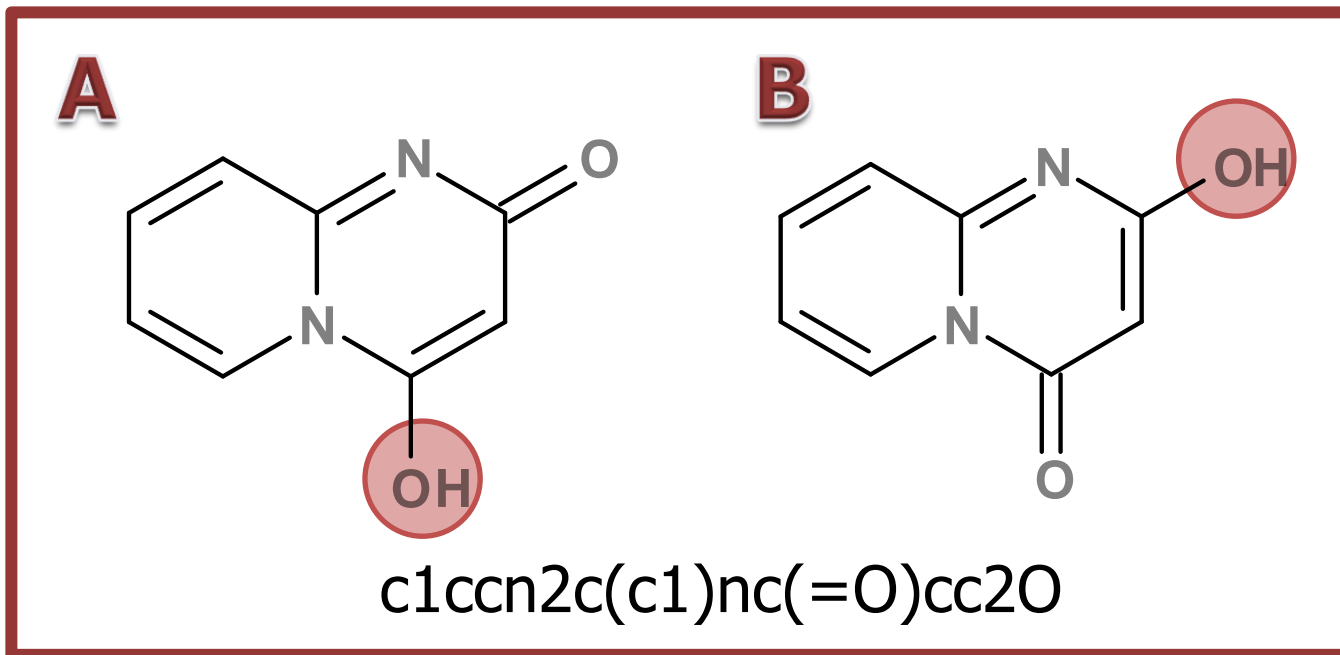


NSC numbers:
98689



NSC numbers:
98690

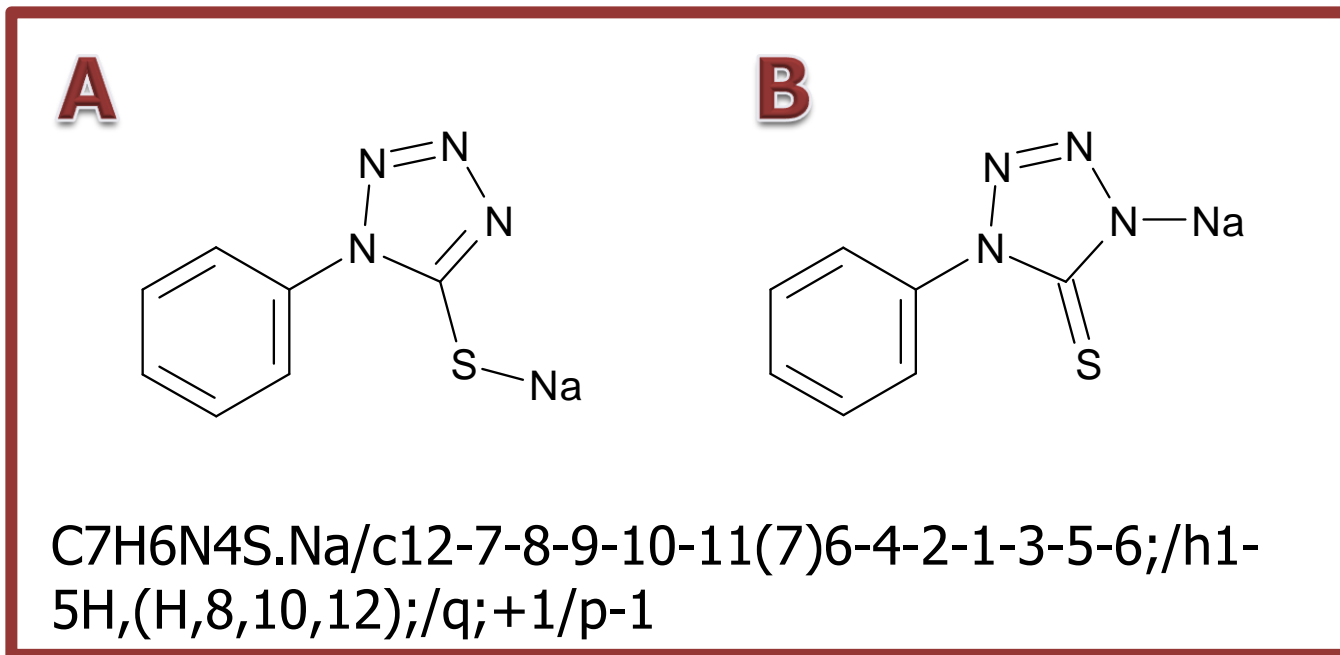
Differences – Tautomer recognition



A C₈H₆N₂O₂/c11-7-5-8(12)10-4-2-1-3-6(10)9-7/h1-5,**12H**

B C₈H₆N₂O₂/c11-7-5-8(12)10-4-2-1-3-6(10)9-7/h1-5,**11H**

InChI Removes Metal Bonds

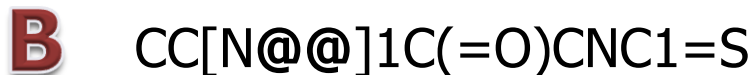
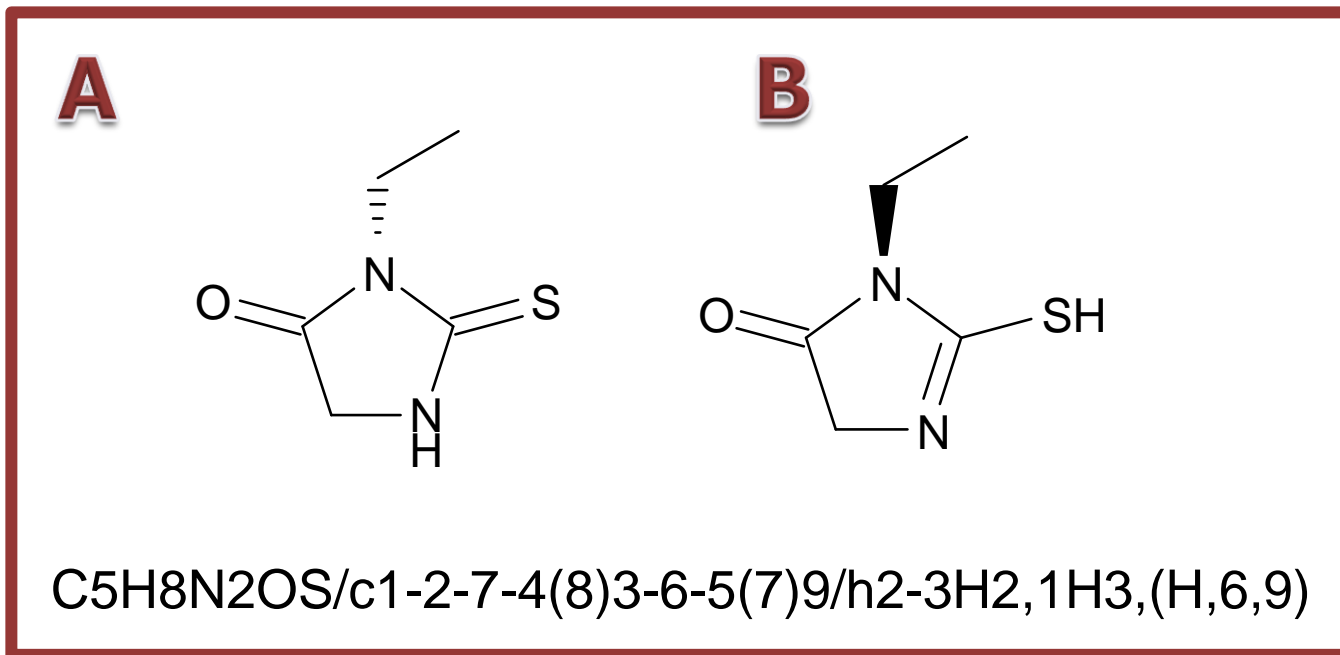


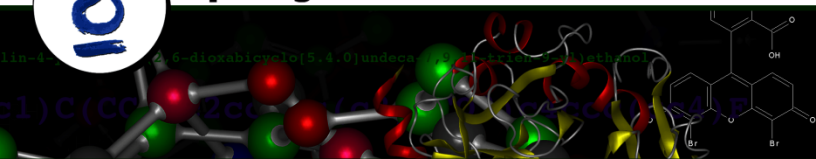
A c1ccc(cc1)[n@]2c(nnn2)S[Na]

B c1ccc(cc1)[n@]2c(=S)[n@](nn2)[Na]



InChI Automatically Cleans Stereo





InChI

- Consistent unique identifier (because reference implementation)
- Normalization/tautomer identification part of the language

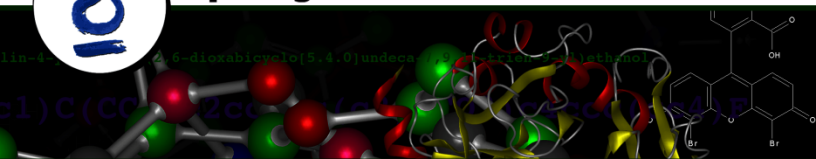
More consistent!

SMILES

- Algorithm dependent unique identifier
- Tautomer handling can be performed by separate application

More flexible!

	INCHI	SMILES
Human readable/writable	NO	YES
Compact	NO	YES
Expresses similarity (LINGOs)	NO	YES



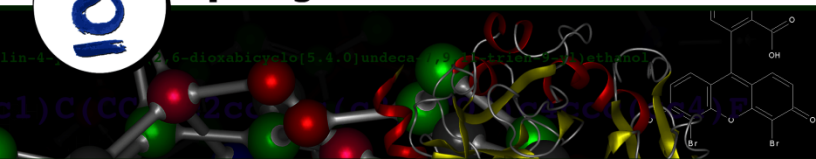
Acknowledgments



Roger Sayle



Brian Cole
Ben Ellingson
Paul Hawkins



Thank you for your attention!

OpenEye Scientific Software

For more information, please contact us.

business@eyesopen.com

support@eyesopen.com

www.eyesopen.com

505-473-7385