

DM840 Algorithms in Cheminformatics

Daniel Merkle

September 5, 2022

Subject overview

- Representation of Molecular Structures
- Combinatorial Structures (Counting, Generating Functions, ...)
- Structure Representation (Canonicalization Algorithms)
- Structure Invariant (e.g. Topological Indices)
- Graph Grammars ("Formal Languages" for Graphs)
- Synthesis Planning (e.g., Shortest Paths in Hypergraphs)
- Enzymatic Design (Discrete Optimization, ILP)
- Concurrency Theory and Causality (e.g. Petri Nets, Category Theory)
- Artificial Chemistries (e.g. "Lattices")
- Quantitative Structure Activity Relationship
 - Principal Component Analysis
 - Algorithms for Minimum Cycle Basis
- Organization Theory
- Stoichiometric Models
- Metabolic Networks and Metabolic Pathways
- (Flux Balance Analysis)
- ...

Representation of Molecular Structures



Figure: From the peyote cactus (*Lophophora williamsii*)

Representation of Molecular Structures

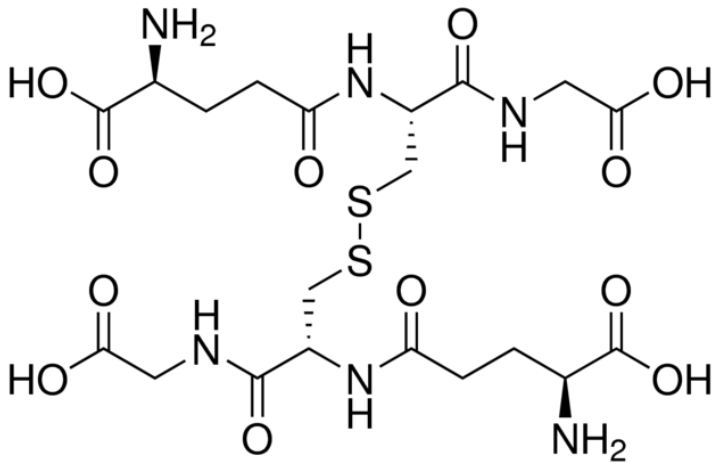


Figure: L-Glutathione oxidized

<https://www.sigmaaldrich.com/catalog/product/sigma/g4376?lang=en>

Representation of Molecular Structures

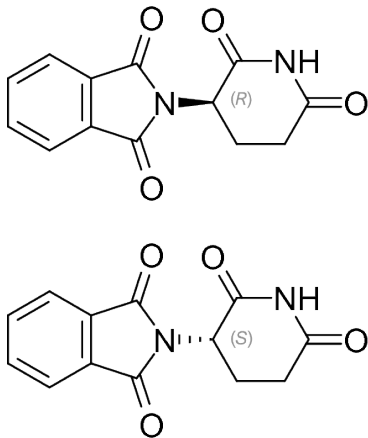
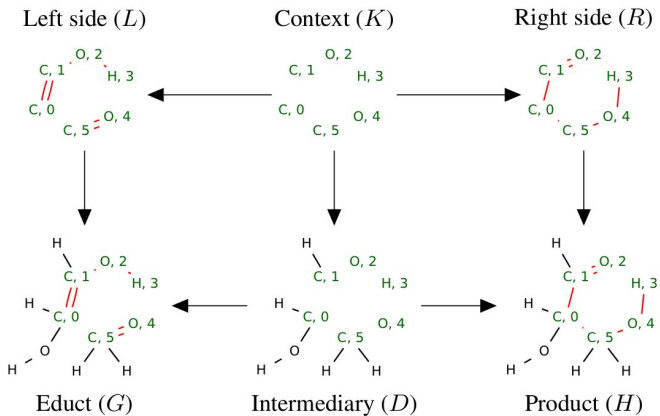
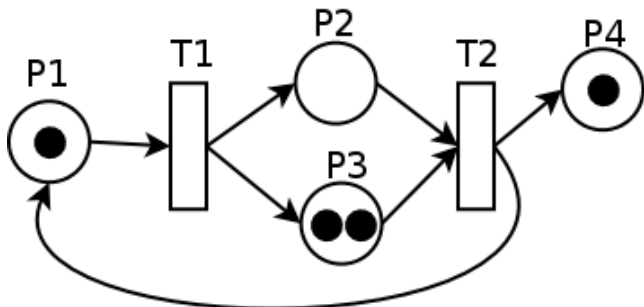


Figure: Thalidomide enantiomers (Contergan)

Graph Grammars

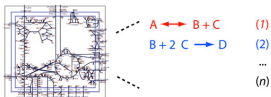


Petri Nets



Flux Balance Analysis

a Curate metabolic reactions



b Formulate **S** matrix

		Reactions			
		1	2	...	n
Metabolites	A	-1			
	B	1	-1		
	C		1	-2	
	D				1
	...				
m					

S

c Apply mass balance constraints

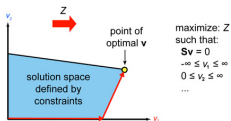
$$\begin{matrix} \mathbf{S} & (m \times n) \\ \begin{bmatrix} 1 & -1 & & \\ 1 & -2 & & \\ & & & \end{bmatrix} \end{matrix} * \begin{matrix} \mathbf{v} & (n \times 1) \\ \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_n \end{bmatrix} \end{matrix} = \mathbf{0} \rightarrow \begin{matrix} m \text{ mass balance} \\ \text{equations} \\ -v_1 + \dots = 0 \\ v_1 - v_2 + \dots = 0 \\ v_1 - 2v_2 + \dots = 0 \\ v_2 + \dots = 0 \\ \dots \end{matrix}$$

d Define objective function **Z**

$$\mathbf{Z} = \begin{matrix} \mathbf{c}' & (1 \times n) \\ \begin{bmatrix} 1 & 0 & \dots & 0 \end{bmatrix} \end{matrix} * \begin{matrix} \mathbf{v} & (n \times 1) \\ \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_n \end{bmatrix} \end{matrix}$$

sets reaction 1 as the objective

e Optimize **Z** using linear programming

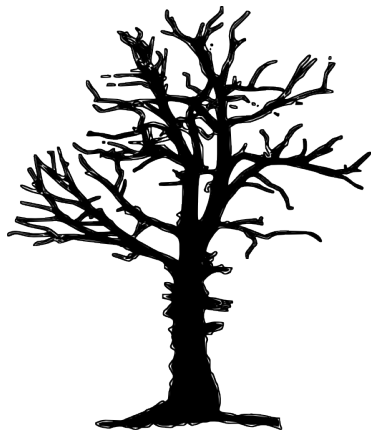


Motivation

Why is the sky blue?



The research situation



Money



Theory - Practice

Theory is when you know everything but nothing works.

Practice is when everything works but no one knows why.

In our lab, theory and practice are combined: nothing works and no one knows why.

The End

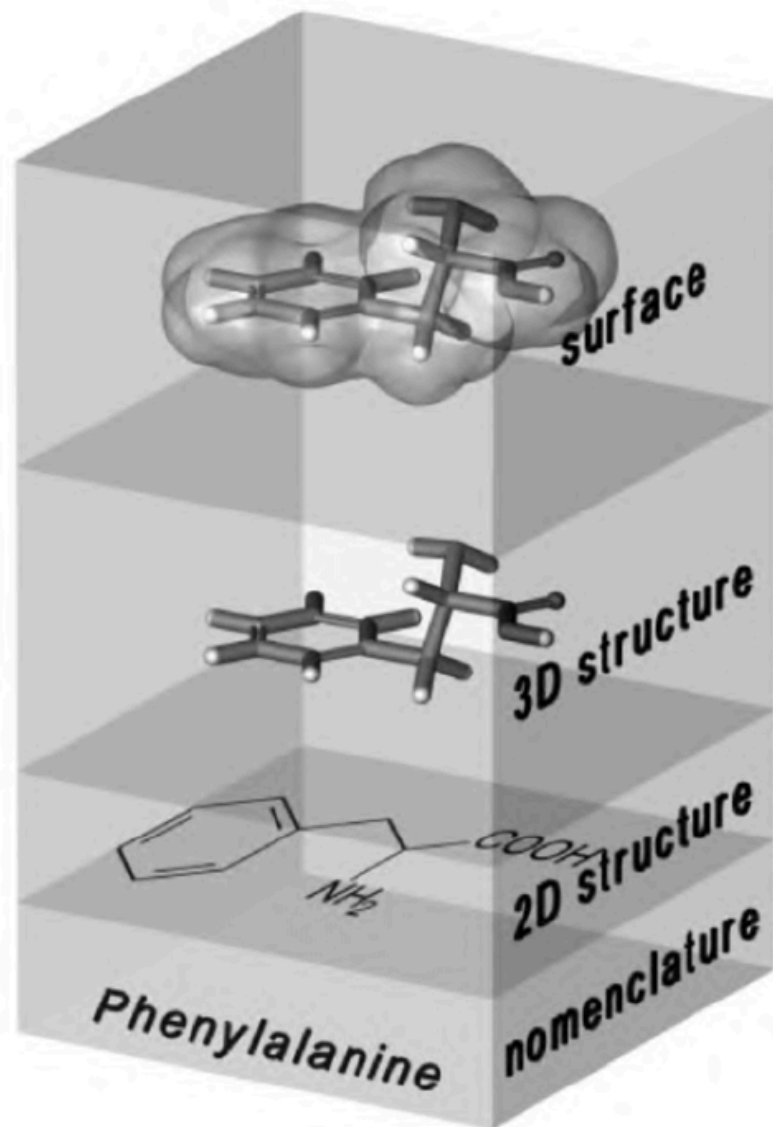
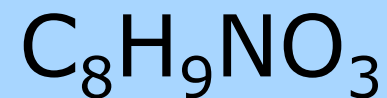


Figure 2-1. Hierarchical scheme for representations of a molecule with different contents of structural information.

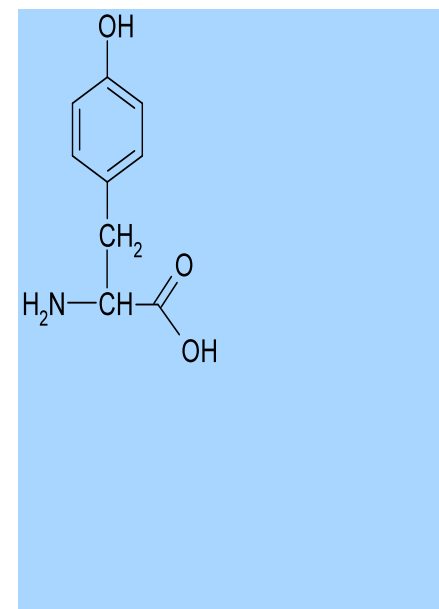
Representing a chemical structure

- How much information do you want to include?
 - atoms present
 - connections between atoms
 - bond types
 - stereochemical configuration
 - charges
 - isotopes
 - 3D-coordinates for atoms



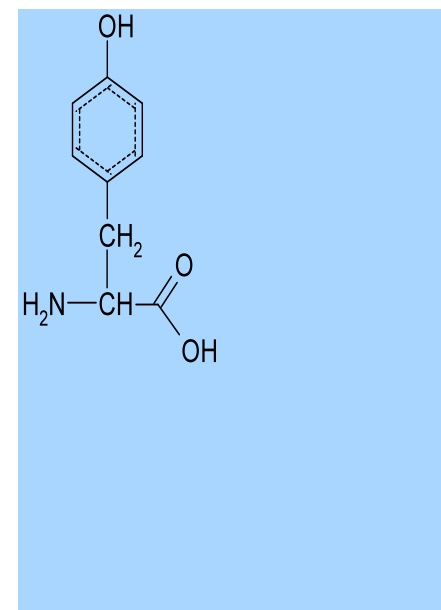
Representing a chemical structure

- How much information do you want to include?
 - atoms present
 - connections between atoms
 - bond types
 - stereochemical configuration
 - charges
 - isotopes
 - 3D-coordinates for atoms



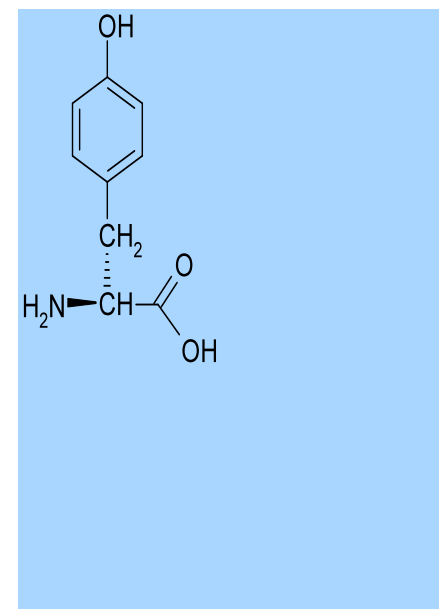
Representing a chemical structure

- How much information do you want to include?
 - atoms present
 - connections between atoms
 - bond types
(aromatic ring identification)
 - stereochemical configuration
 - charges
 - isotopes
 - 3D-coordinates for atoms



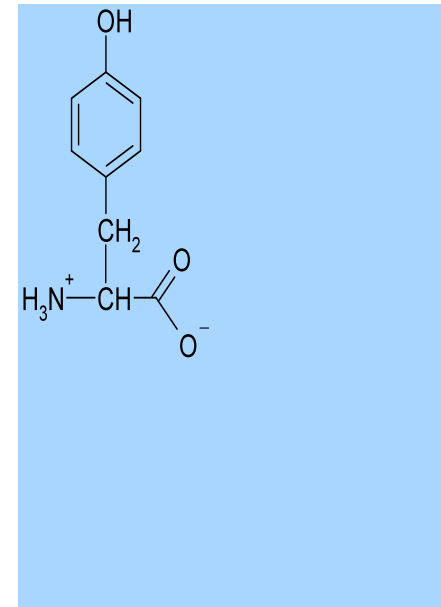
Representing a chemical structure

- How much information do you want to include?
 - atoms present
 - connections between atoms
 - bond types
 - stereochemical configuration
 - charges
 - isotopes
 - 3D-coordinates for atoms



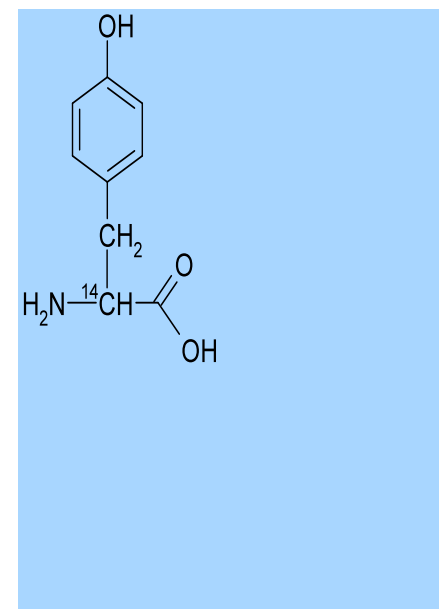
Representing a chemical structure

- How much information do you want to include?
 - atoms present
 - connections between atoms
 - bond types
 - stereochemical configuration
 - charges
 - isotopes
 - 3D-coordinates for atoms



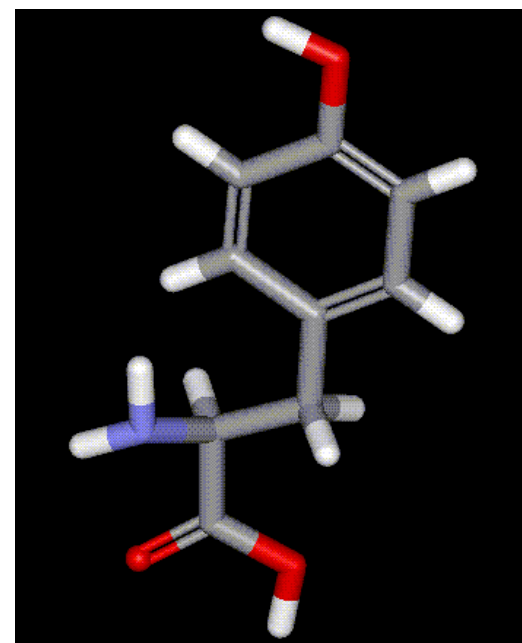
Representing a chemical structure

- How much information do you want to include?
 - atoms present
 - connections between atoms
 - bond types
 - stereochemical configuration
 - charges
 - isotopes
 - 3D-coordinates for atoms

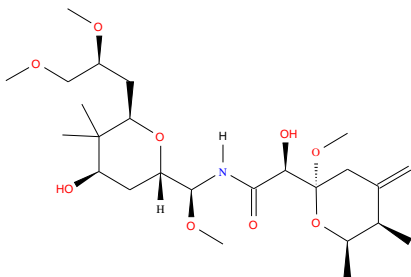


Representing a chemical structure

- How much information do you want to include?
 - atoms present
 - connections between atoms
 - bond types
 - stereochemical configuration
 - charges
 - isotopes
 - 3D-coordinates for atoms



Store Substance as Graphics or by Name



- Pederin
- **2H-Pyran-2-glycolamide**, N-((6-(2,3-dimethoxypropyl) tetrahydro-4-hydroxy-5,5-dimethyl-2H-pyran-2-yl)methoxymethyl) tetrahydro-2-methoxy-5,6-dimethyl-4- methylene-
- **D-manno-Nonitol**, 2,6-anhydro-3,5,7-trideoxy-1-C-(((2S)-hydroxy((2R,5R,6R)-tetrahydro-2-methoxy-5,6-dimethyl-4-methylene-2H-pyran-2-yl)acetyl)amino)-5,5-dimethyl-1,8,9-tri-O-methyl-, (1S)-

Elements of Formal Grammars

- 1 **Terminal Symbols** T (represented by lowercase letters).
- 2 **Nonterminals Symbols** N (represented by uppercase letters).
- 3 **Production Rules** with a left- and a right-hand side consisting of strings of these symbols.
- 4 **Start Symbol** (also called Axiom)

The example grammar defines the language of all strings of the form $\{ax_1x_2\dots x_kb \mid k \geq 0 \wedge x_i \in \{a, b\}\}$. (A is the Axiom).

members: ab, abab, aaabbbb

non-members: a, b, ababa

$$T = \{a, b\}$$

$$N = \{A, B, C\}$$

$$A \rightarrow aA \mid aB$$

$$B \rightarrow aB \mid bB \mid C$$

$$C \rightarrow bC \mid b$$

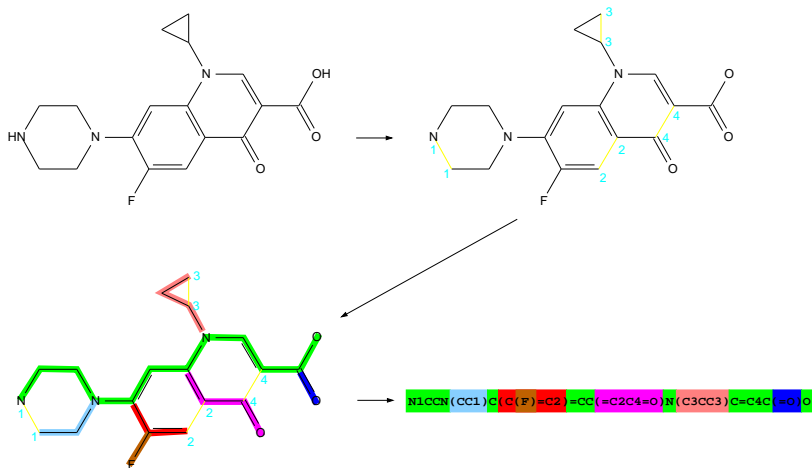
Chomsky Hierarchy

Type No	Type Name	Rule Pattern	
0	unrestricted	$x \rightarrow y,$	$x \in (N \cup T)_+$ $y \in (N \cup T)^*$
1	context sensitive	$x \rightarrow y,$	$x \in (N \cup T)_+$ $y \in (N \cup T)_+$ $ x \leq y $
2	context free	$x \rightarrow y,$	$x \in N$ $y \in (N \cup T)^*$
3	regular	$w \rightarrow x \mid yz$	$w, x, z \in N$ $y \in T$

BNF grammar of Daylight's SMILES

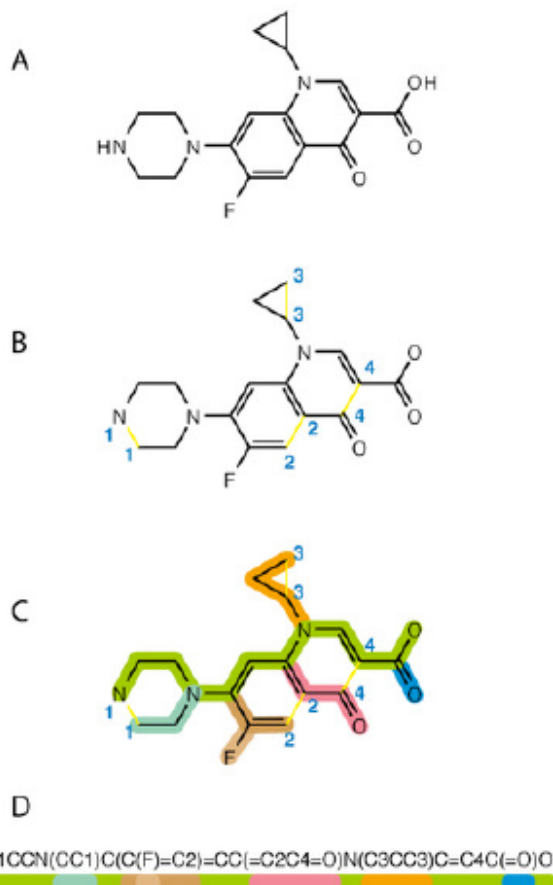
```
smiles          ::= chain terminator
chain           ::= branched_atom | chain branched_atom
                 | chain bond branched_atom | chain '.' branched_atom
branched_atom  ::= atom ringbond* branch*
atom           ::= bracket_atom | aliphatic_organic | aromatic_organic
                 | '*'
ringbond       ::= bond? DIGIT | bond? '\%' DIGIT DIGIT
branch         ::= '(' chain ')' | '(' bond chain ')' | '(' '.' chain ')'
bracket_atom   ::= '[' isotope? symbol chiral? hcount? charge? class? ']'
isotope        ::= NUMBER
symbol         ::= element_symbols | aromatic_symbols | '*'
chiral         ::= '@' | '@@'
hcount         ::= 'H' DIGIT?
charge         ::= '-' DIGIT? | '+' DIGIT?
class         ::= ':' NUMBER
aliphatic_organic ::= 'B' | 'C' | 'N' | 'O' | 'S' | 'P' | 'F' | 'Cl' | 'Br' | 'I'
aromatic_organic ::= 'b' | 'c' | 'n' | 'o' | 's' | 'p'
element_symbols ::= 'H' | 'He' | 'Li' | 'Be' | 'B' | 'C' | 'N' | 'O' | etc
aromatic_symbols ::= 'c' | 'n' | 'o' | 'p' | 's' | 'se' | 'as'
bond           ::= '-' | '=' | '#' | '\$' | ':' | '/' | '\'
terminator     ::= SPACE | TAB | '\n' | '\0'
```

Generation of SMILES



modified from <http://en.wikipedia.org/wiki/File:SMILES.png>

SMILES (Simplified Molecular Input Line Entry System)



Six basic rules:



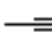

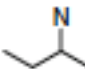
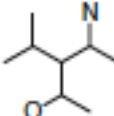
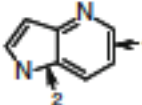
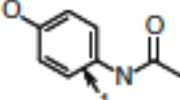
- 1) Atoms by atomic symbol
- 2) Hydrogen atoms added to saturate free valences
- 3) Neighboring atoms stand next to one another
- 4) Double and triple bonds by = and #
- 5) Branching shown by parentheses
- 6) Rings shown by digit at ring closures

Canonical SMILES: unique for each structure

Isomeric SMILES: describe isotopism, configuration around double bonds and tetrahedral centers, chirality

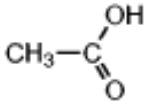


SMILES

Illustrative SMILES: molecular structures and the corresponding SMILES strings are paired vertically. The numbered arrows on the three cyclic molecular structures are not part of the molecules. They are used to indicate the break points for deriving the corresponding SMILES strings (see text)

			
<chem>CCC</chem>	<chem>CC=C</chem>	<chem>CC#C</chem>	<chem>c1ccncc1</chem>
			
<chem>CCC(C)N</chem>	<chem>CC(C)C(C(C)N)C(C)O</chem>	<chem>c1cc2c(cc[nH]2)nc1</chem>	<chem>CC(-O)Nc1ccc(cc1)O</chem>

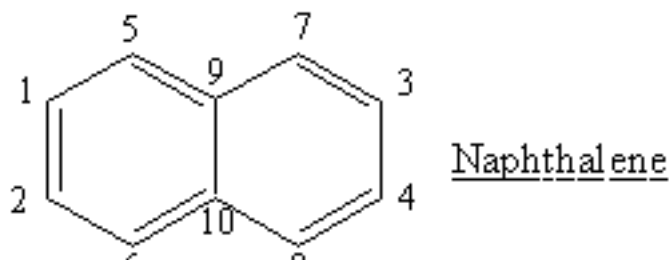
SLN (SYBYL Line Notation)

Table 2-3. Basic SLN syntax without description of attributes and macro atoms.

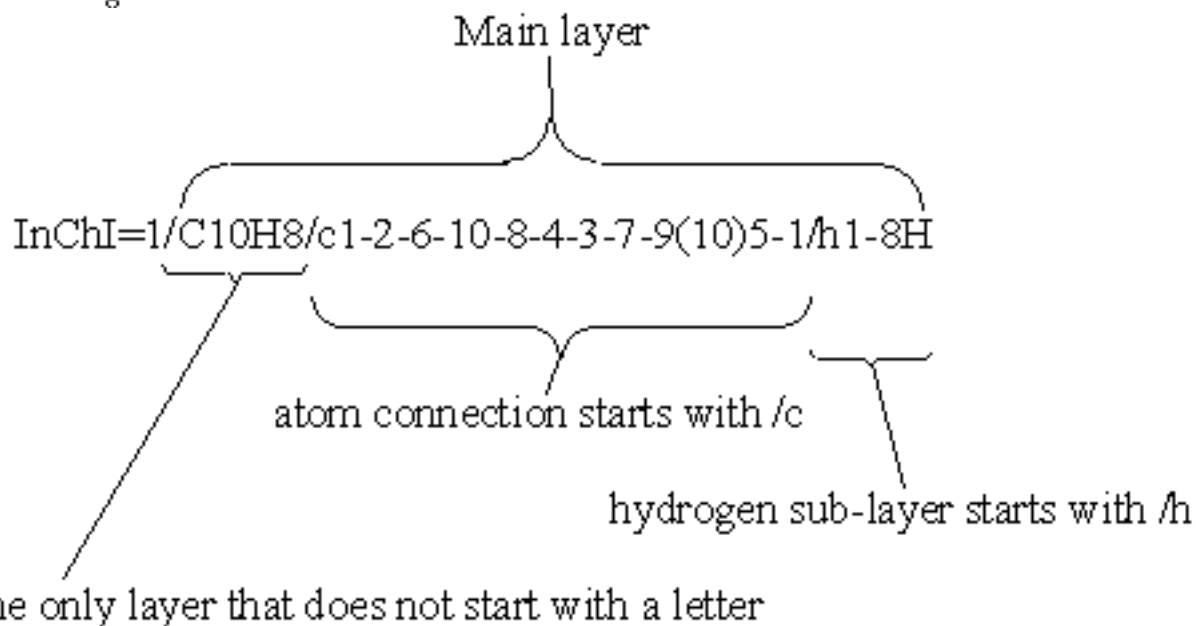
<i>SLN</i>	<i>Chemical structure</i>	<i>Compound name</i>
<i>Atoms:</i> Atoms are represented by their atomic symbols. The first letter is upper-case, and in two-letter symbols the second letter is lower-case. Hydrogen atoms must be specified.		
CH4	CH ₄	methane
NH2	-NH ₂	amine
<i>Bonds:</i> Single bonds are omitted; double, triple, and aromatic bonds are indicated by the symbols "=", "#", and ":", respectively. In contrast to SMILES, aromaticity is not an atomic property, but a property of bonds. A period indicates the start of a new part of the structure.		
HC(=O)OH	HCOOH	formic acid
Na.OH	NaOH	sodium hydroxide
<i>Branches:</i> Branches are indicated by parentheses.		
CH3C(=O)OH		acetic acid
<i>Cyclic structures:</i> Ring closures are described by a bond to a previously defined atom which is specified by a unique ID number. The ID is a positive integer placed in square brackets behind the atom. An "@" indicates a ring closure.		
C[15]H2CH2CH2CH2CH2@15		cyclohexane
O[6]:CH:CH:CH:CH:@6		furan

InChI (IUPAC International Chemical Identifier)

InChI=1/C10H8/c1-2-6-10-8-4-3-7-9(10)5-1/h1-8H



1. Main layer
2. Charge layer
3. Stereochemical layer
4. Isotopic layer
5. Fixed-H layer
6. Reconnected Layer



Adjacency and Distance matrices

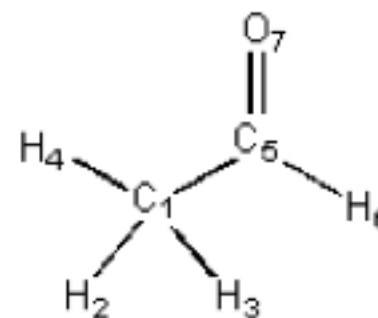
Acetaldehyde: CH₃CH=O

Adjacency Matrix

C1	0	1	1	1	1	0	0
H2	1	0	0	0	0	0	0
H3	1	0	0	0	0	0	0
H4	1	0	0	0	0	0	0
C5	1	0	0	0	0	1	1
H6	0	0	0	0	1	0	0
O7	0	0	0	0	1	0	0

Distance Matrix

0	1	1	1	1	2	2
1	0	2	2	2	3	3
1	2	0	2	2	3	3
1	2	2	0	2	3	3
1	2	2	2	0	1	1
2	3	3	3	1	0	2
2	3	3	3	1	2	0



1: atoms i j are bonded

0: atoms i j are not bonded

Length of shortest path
between atoms i j

Bond matrix

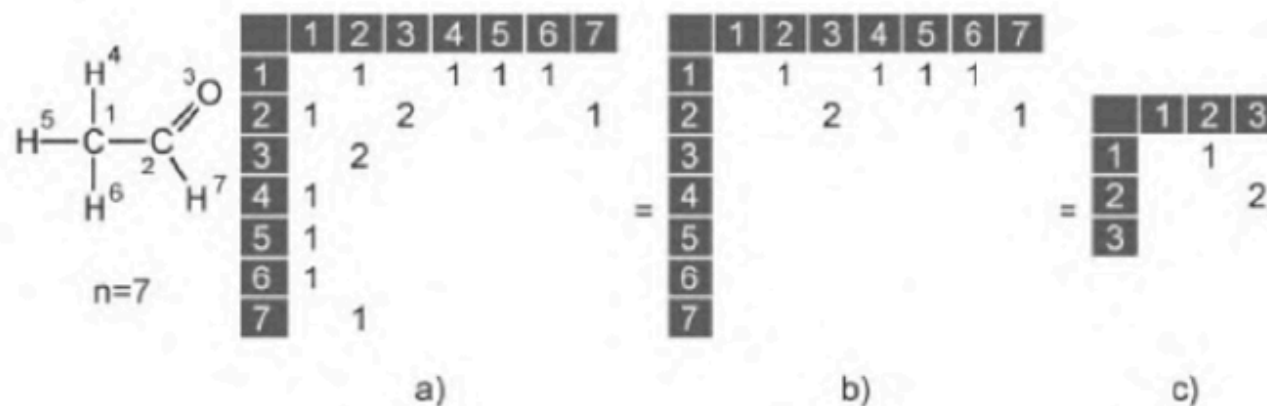
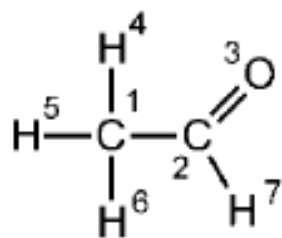


Figure 2-17. a) The redundant bond matrix of ethanal with the zero values omitted. b) It can be compressed by reduction to the top right triangle. c) Omitting the hydrogen atoms provides the simplest non-redundant matrix representation.

Connection Table



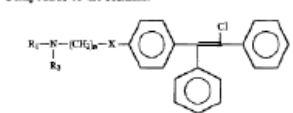
Atom list	
1	C
2	C
3	O
4	H
5	H
6	H
7	H

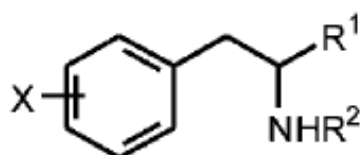
Bond list		
1 st atom	2 nd atom	bond order
1	2	1
2	3	2
2	7	1
1	4	1
1	5	1
1	6	1

Figure 2-20. A connection table: the structure diagram of ethanal, with the atoms arbitrarily labeled, is defined by a list of atoms and a list of bonds.

Special notations of chemical structures

- Markush structures

United States Patent [19]		[11] Patent Number:	5,681,863
Bitonti et al.		[45] Date of Patent:	Oct. 28, 1997
[54]	NON-METABOLIZABLE CLOMPHENE ANALOGS FOR TREATMENT OF TAMOXIFEN-RESISTANT TUMORS	3,631,109 12/1971 O'Sega et al.	269370.9
		4,696,949 9/1987 Tavola et al.	514664
		4,839,155 6/1989 McCague	514651
		5,114,951 9/1992 King et al.	514950
		5,130,424 7/1992 Waintrub	54028
[75]	Inventors: Alan J. Bitonti, Malneville; Russell J. Bammann, Cincinnati, both of Ohio	<i>Primary Examiner</i> —Jerome D. Goldberg <i>Attorney Agent, or Firm</i> —Nelson L. Lentz	
[73]	Assignee: Merrell Pharmaceuticals Inc., Cincinnati, Ohio	[57] ABSTRACT	
[21]	Appl. No.: 350,192	Compounds of the formula:	
[22]	Filed: Dec. 5, 1994		
Related U.S. Application Data			
[62]	Division of Ser. No. 196,817, Feb. 10, 1994, Pat. No. 5,410,010, which is a continuation of Ser. No. 945,305, Sep. 15, 1992, abandoned.	wherein R ₁ and R ₂ are each selected from the group consisting of C ₁ -C ₄ lower alkyl; X is NH or S; and n is a whole number within the range of 1-4 inclusive; and when n=0, X is (CH ₂) ₂ and the pharmaceutically acceptable salts thereof	
[51]	Int. Cl. ⁵	A61K 31/135	
[52]	U.S. Cl.	514/648; 514/649	
[58]	Field of Search	514/641, 649	
[56]	References Cited		



R¹ = H or small alkyl, halogen, OH, COOH

R² = H, CH₃

X = H, (CH₂)_nCH₃

Figure 2-62. The substituted phenyl derivative is an example of a typical Markush structure. Herein, a number of compounds are described in one structure diagram by fill-ins. Phenylalanine is one of these structures when R¹ is COOH, R² is H, and X is H.

Special notations of chemical structures

- Fingerprints

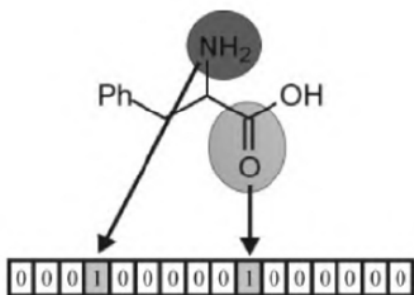


Figure 2-64. How an excerpt from a binary code could appear, if only -NH_2 and $\text{C}=\text{O}$ are available in the fragment library.

MACCS fingerprints: 166 structural keys

that answer questions of the type:

- Is there a ring of size 4?
- Is at least one F, Br, Cl, or I present?

where the answer is either

TRUE (1) or FALSE (0)

SMILES .smi file

```
N12CCC36C1CC(C(C2)=CCOC4CC5=O)C4C3N5c7ccccc76 Strychnine
c1ccccc1C(=O)OC2CC(N3C)CCC3C2C(=O)OC cocaine
COc1cc2c(ccnc2cc1)C(O)C4CC(CC3)C(C=C)CN34 quinine
OC(=O)C1CN(C)C2CC3=CCNc(ccc4)c3c4C2=C1 lyseric acid
CCN(CC)C(=O)C1CN(C)C2CC3=CNc(ccc4)c3c4C2=C1 LSD
C123C5C(O)C=CC2C(N(C)CC1)Cc(ccc4O)c3c4O5 morphine
C123C5C(OC(=O)C)C=CC2C(N(C)CC1)Cc(ccc4OC(=O)C)c3c4O5 heroin
c1ncccc1C1CCCN1C nicotine
CN1C(=O)N(C)C(=O)C(N(C)C=N2)=C12 caffeine
C1C(C)=C(C=CC(C)=CC=CC(C)=CCO)C(C)(C)C1 vitamin a
```

MOLfile

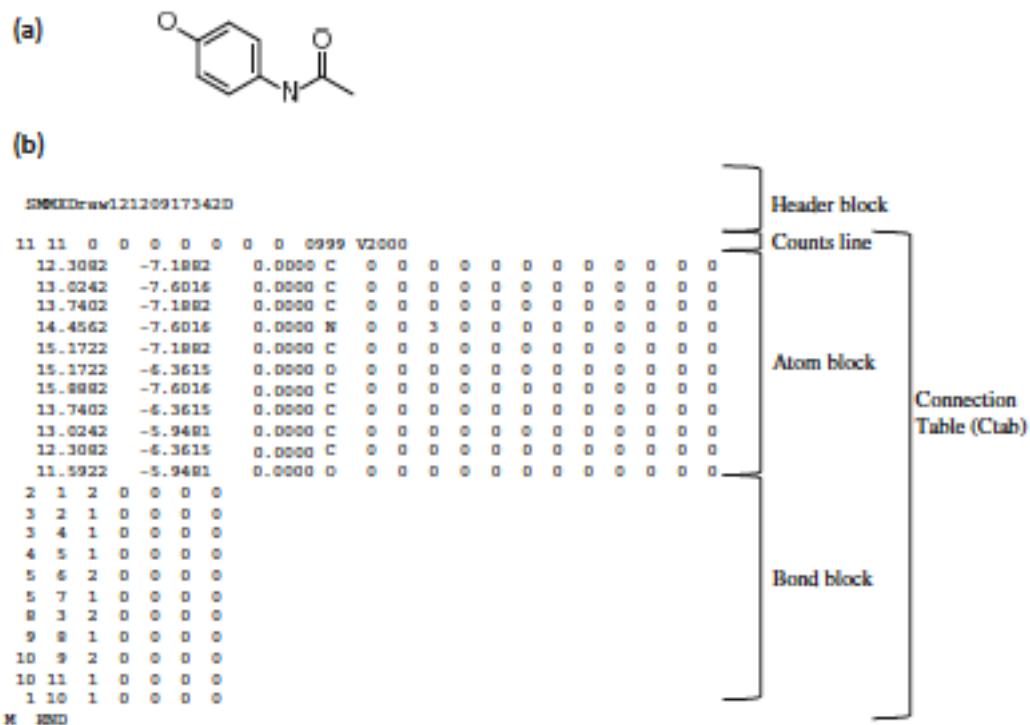


Fig. 2.1. Illustrative example of a MOLfile for acetaminophen (also known as paracetamol). (a) Molecular structure of acetaminophen, commonly known as Tylenol. Tylenol is a widely used medicine for reducing fever and pain. (b) MOLfile for acetaminophen.

Aspirin in SDF Format

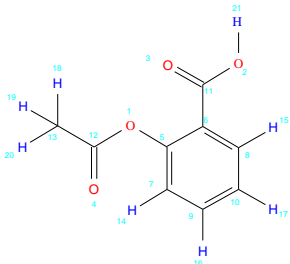
2244

-OEChem-03090904423D

```
21 21 0      0 0 0 0 0 0 0999 V2000
 1.8152 -0.9382  4.0419 O  0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
 5.1920 -2.1043  2.0467 O  0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
 3.9623 -2.6855  3.8563 O  0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
 2.9441  1.1113  3.9712 O  0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
 1.8509 -0.9767  2.6799 C  0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
 2.9180 -1.5734  2.0082 C  0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
 0.8008 -0.4105  1.9570 C  0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
 2.9348 -1.6038  0.6137 C  0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
 0.8176 -0.4410  0.5626 C  0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
 1.8846 -1.0376 -0.1090 C  0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
 4.0236 -2.1714  2.7435 C  0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
 2.4171  0.2017  4.5978 C  0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
 2.3265  0.1503  6.0942 C  0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
-0.0345  0.0550  2.4732 H  0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
 3.7445 -2.0729  0.0609 H  0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
-0.0005 -0.0011 -0.0002 H  0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
 1.8958 -1.0640 -1.1947 H  0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
 1.2777  0.1264  6.3998 H  0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
 2.7936  1.0443  6.5170 H  0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
 2.8566 -0.7302  6.4654 H  0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
 5.9361 -2.5075  2.5428 H  0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0

 1  5  1  0  0  0  0  0
 1 12  1  0  0  0  0  0
 2 11  1  0  0  0  0  0
 2 21  1  0  0  0  0  0
 3 11  2  0  0  0  0  0
 4 12  2  0  0  0  0  0
 5  6  1  0  0  0  0  0
 5  7  2  0  0  0  0  0
 6  8  2  0  0  0  0  0
 6 11  1  0  0  0  0  0
 7  9  1  0  0  0  0  0
 7 14  1  0  0  0  0  0
 8 10  1  0  0  0  0  0
 8 15  1  0  0  0  0  0
 9 10  2  0  0  0  0  0
 9 16  1  0  0  0  0  0
10 17  1  0  0  0  0  0
12 13  1  0  0  0  0  0
13 18  1  0  0  0  0  0
13 19  1  0  0  0  0  0
13 20  1  0  0  0  0  0
```

M END



2244

-OEChem-03090904423D

21	21	0	0	0	0	0	0	0999	v2000			
	1.8152		-0.9382			4.0419	O		0	0	0	
	5.1920		-2.1043			2.0467	O		0	0	0	
	3.9623		-2.6855			3.8563	O		0	0	0	
	2.9441		1.1113			3.9712	O		0	0	0	
	1.8509		-0.9767			2.6799	C		0	0	0	
	2.9180		-1.5734			2.0082	C		0	0	0	
	0.8008		-0.4105			1.9570	C		0	0	0	
	2.9348		-1.6038			0.6137	C		0	0	0	
	0.8176		-0.4410			0.5626	C		0	0	0	
	1.8846		-1.0376			-0.1090	C		0	0	0	
	4.0236		-2.1714			2.7435	C		0	0	0	
	2.4171		0.2017			4.5978	C		0	0	0	
	2.3265		0.1503			6.0942	C		0	0	0	

2.8566

-0.7302

6.4654

5.9361

-2.5075

2.5428

1 5 1 0 0 0 0

1 12 1 0 0 0 0

2 11 1 0 0 0 0

2 21 1 0 0 0 0

3 11 2 0 0 0 0

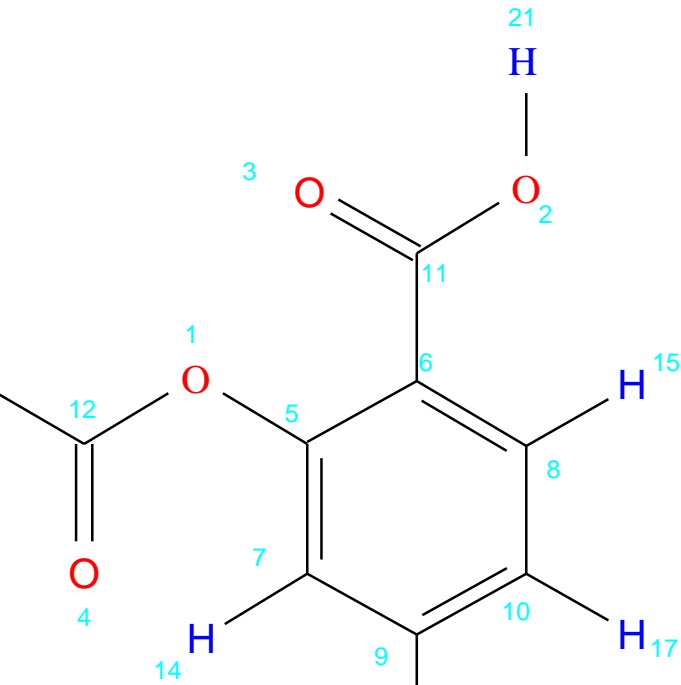
4 12 2 0 0 0 0

5 6 1 0 0 0 0

5 7 2 0 0 0 0

6 8 2 0 0 0 0

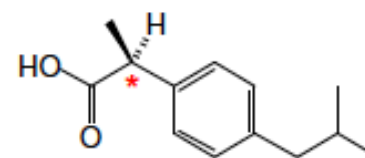
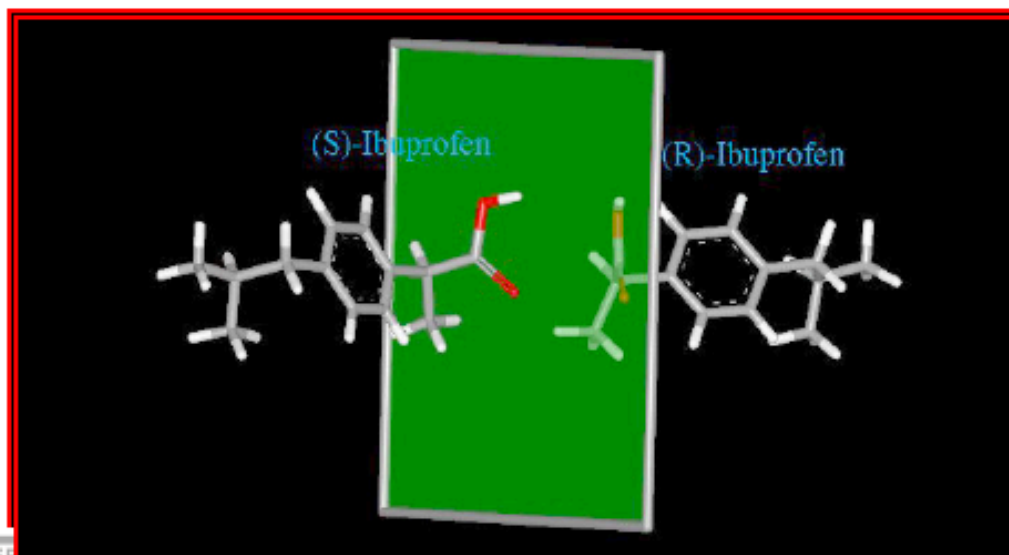
6 11 1 0 0 0 0



Importance of stereochemistry

Enantiomers (mirror image molecules) have

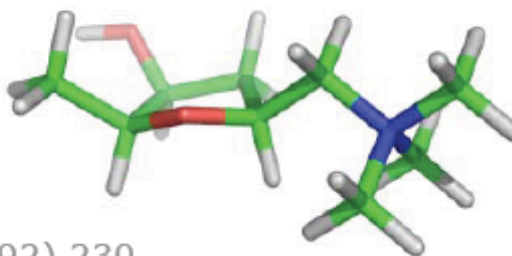
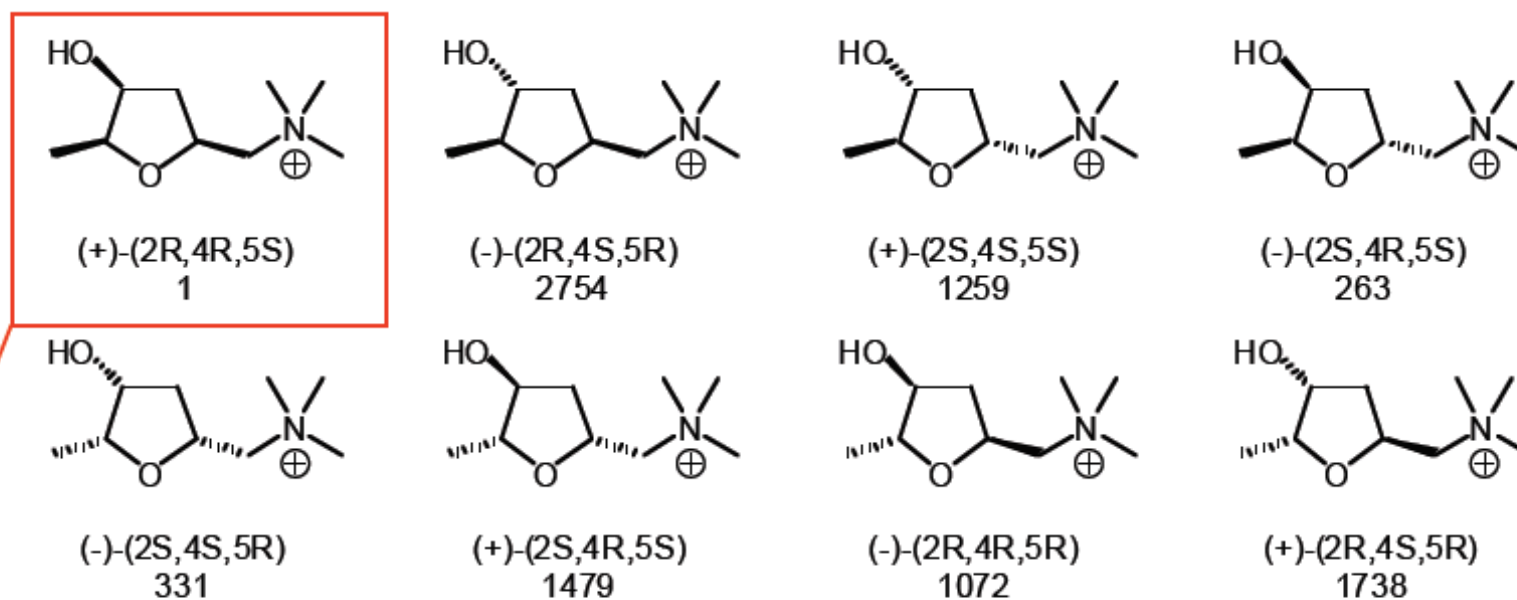
- ☆ identical physical and chemical properties
- ☆ different optical rotation
- ☆ often different biological activities



3D link

Importance of stereochemistry

- Muscarine (muscarinic agonist): Natural product most potent
- 7 other stereoisomers possible: Considerable less potent



isomers

compounds have an identical empirical formula, but different molecular structures

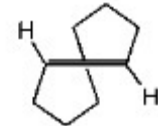

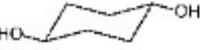
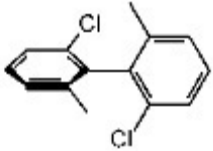

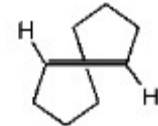
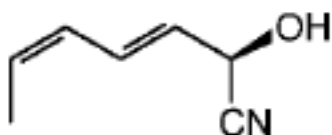
constitutional isomers		stereoisomers			
structural isomers	positional isomers	configurational isomers		conformational isomers	
<p>different functional groups</p> <div style="display: flex; justify-content: space-around; align-items: center;"> <div style="text-align: center;"> $\begin{array}{c} \text{CH}_3 \\ \\ \text{O} \\ \\ \text{CH}_3 \end{array}$ dimethyl ether </div> <div style="text-align: center;"> $\begin{array}{c} \text{CH}_3 \\ \\ \text{CH}_2 \\ \\ \text{OH} \end{array}$ ethanol </div> </div>	<p>identical functional groups at different places</p> <div style="display: flex; justify-content: space-around; align-items: center;"> <div style="text-align: center;"> $\begin{array}{c} \text{COOH} \\ \\ \text{H}_2\text{N}-\text{C}-\text{H} \\ \\ \text{CH}_3 \end{array}$ α-alanine </div> <div style="text-align: center;"> $\begin{array}{c} \text{COOH} \\ \\ \text{CH}_2 \\ \\ \text{H}_2\text{N}-\text{CH}_2 \end{array}$ β-alanine </div> </div>	chirality			<p>different conformers by rotation around a single bond</p> <div style="display: flex; justify-content: space-around; align-items: center;"> <div style="text-align: center;"> $\begin{array}{c} \text{CH}_3 \\ \\ \text{H}-\text{C}-\text{C}-\text{H} \\ \quad \\ \text{H} \quad \text{H} \end{array}$ <i>gauche</i> </div> <div style="text-align: center;"> $\begin{array}{c} \text{CH}_3 \\ \\ \text{H}-\text{C}-\text{C}-\text{H} \\ \quad \\ \text{H} \quad \text{CH}_3 \end{array}$ <i>n-butane anti</i> </div> </div> <div style="text-align: center; margin-top: 10px;">  E-cyclooctene </div>
		cis/trans isomers	chiral center		
<p>neighboring groups of double bonds can have two directions.</p> <div style="display: flex; justify-content: space-around; align-items: center;"> <div style="text-align: center;"> $\begin{array}{c} \text{H} \quad \text{COOH} \\ \diagdown \quad / \\ \text{C}=\text{C} \\ / \quad \diagdown \\ \text{H} \quad \text{COOH} \end{array}$ maleic acid (<i>cis</i> = <i>Z</i>) </div> <div style="text-align: center;"> $\begin{array}{c} \text{H} \quad \text{COOH} \\ \diagdown \quad / \\ \text{C}=\text{C} \\ / \quad \diagdown \\ \text{HOOC} \quad \text{H} \end{array}$ fumaric acid (<i>trans</i> = <i>E</i>) </div> </div> <p>Ring systems can also be differentiated into <i>cis/trans</i> arrangements of substituents.</p> <div style="display: flex; justify-content: space-around; align-items: center;"> <div style="text-align: center;">  <i>cis</i>-1,4-cyclohexanediol </div> <div style="text-align: center;">  <i>trans</i>- </div> </div>	<p>if a molecule has one chiral atom, the isomers are like an image and mirror image (enantiomers); molecules with more than one chiral atom exist as enantiomers and diastereomers</p> <div style="display: flex; justify-content: space-around; align-items: center;"> <div style="text-align: center;"> $\begin{array}{c} \text{COOH} \\ \\ \text{H}-\text{C}-\text{NH}_2 \\ \\ \text{CH}_3 \end{array}$ D-alanine </div> <div style="text-align: center;"> $\begin{array}{c} \text{COOH} \\ \\ \text{H}_2\text{N}-\text{C}-\text{H} \\ \\ \text{CH}_3 \end{array}$ L-alanine </div> </div>	<p>if a molecule has four ligands which are placed pairwise along an axis and are not in one plane (atropisomer)</p> <div style="text-align: center; margin: 10px 0;">  2,2'-dichloro-6,6'-dimethyl-1,1'-biphenyl </div> <p>a special case of axial chirality is the arrangement of the molecule as a right- or left-handed helix (helicene)</p> <div style="text-align: center; margin-top: 10px;">  heptahelicene </div>	<p>if the arrangement of a molecule can be distinguished into different sides</p> <div style="text-align: center; margin-top: 10px;">  E-cyclooctene </div>		

Figure 2-67. Classification of isomeric structures of organic compounds.

Stereochemistry in SMILES



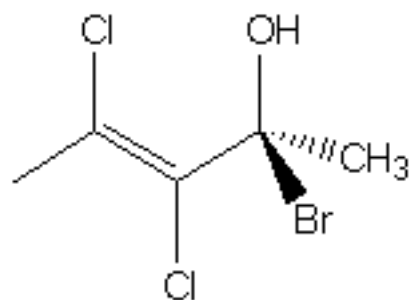
SMILES: C/C=C\C=[C@@H](O)C#N

Figure 2-78. The stereochemistry of (2*R*,3*E*,5*E*)-2-hydroxy-3,5-heptadiene nitrile can be expressed in the SMILES notation with @ or (back)slashes.

@@: when viewed from atom along the bond to the chiral center, the sequence of atoms (H), (O) and (C#N) appear clockwise

"/" and "\": "cis" and "trans" configuration

Stereochemistry in InChI



InChI=1/C5H7BrCl2O/c1-3(7)4(8)5(2,6)9/h9H,1-2H3/b4-3+/t5-/m1/s1

Stereochemical layer

Double bond stereo sub-layer

sp³ stereo sub-layer

File formats

Table 2-5. The most important file formats for exchange of chemical structure information.

<i>File format</i>	<i>Suffix</i>	<i>Comments</i>	<i>Support</i>	<i>Ref.</i>
MDL Molfile	*.mol	Molfile; the most widely used connection table format	www.mdli.com	50
SDfile	*.sdf	Structure-Data file; extension of the MDL Molfile containing one or more compounds	www.mdli.com	50
RDfile	*.rdf	Reaction-Data file; extension of the MDL Molfile containing one or more sets of reactions	www.mdli.com	50
SMILES	*.smi	SMILES; the most widely used linear code and file format	www.daylight.com	20, 21
PDB file	*.pdb	Protein Data Bank file; format for 3D structure information on proteins and polynucleotides	www.rcsb.org	53
CIF	*.cif	Crystallographic Information File format; for 3D structure information on organic molecules	www.iucr.org/iucr-top/cif/	55
JCAMP	*.jdx, *.dx, *.cs	Joint Committee on Atomic and Molecular Physical Data; structure and spectroscopic format	www.jcamp.org/	56
CML	*.cml	Chemical Markup Language; extension of XML with specialization in chemistry	www.xml-cml.org	57–59

Molecular editors and viewers

<http://www.chemaxon.com/products/marvin/>

The screenshot displays the ChemAxon website interface. At the top left is the ChemAxon logo. A navigation bar includes links for Home, Products, Download, Documentation, Forum, Support, About us, Webshop, and Contact us. Below this is a breadcrumb trail: Home > Products > Marvinsketch > Marvinsketch. A search bar is located on the right side of the navigation bar. The main content area features a large heading "Advanced chemical drawing software" and a sub-heading "Marvinsketch is an advanced chemical editor for drawing chemical structures, queries and reactions". A prominent blue button with a download icon says "Download Marvinsketch Suite Sketch/Space/View - Ver 5.5.1.0". To the right of this button is a grid of various chemical structures. Below the main heading is a section titled "Easy chemical structure, query and reaction drawing", which is divided into three columns: "Chemical structure drawing", "Reaction drawing", and "Query drawing". Each column contains a brief description of the software's capabilities in that area. On the left side of the page, there is a vertical menu with categories like "Products", "Naming", "Workflow tools", and "Try now". The "Try now" section lists "Marvinsketch", "MarvinView", and "Calculator Plugins". On the right side, there is a "Related links" section with links to "Release Notes", "History of changes", "User's Guide", "Developer's Guide", "Example implementations", "API", "Feature animations", and "Articles and Presentations".

ChemAxon

Home > Products > Marvinsketch > Marvinsketch

Search

Products >>

- Marvin >>
- Marvinsketch**
- MarvinView
- MarvinSpace
- MolConverter

Calculator Plugins >>

- Naming
- JChem for Excel
- Instant JChem
- JChem Base
- JChem Cartridge
- JChem Web Services Add-on

Workflow tools

- JChem for SharePoint
- Standardizer
- Structure Checker
- Markush Search Add-on
- Screen
- JKluster
- Reactor
- Metabolizer Preview
- Fragmenter
- Online tryouts >>

Try now

- Marvinsketch
- MarvinView
- Calculator Plugins

Advanced chemical drawing software

Marvinsketch is an advanced chemical editor for drawing chemical structures, queries and reactions

Download Marvinsketch Suite
Sketch/Space/View - Ver 5.5.1.0

Try online!

- PDF Brochure
- PPT Technical presentation
- Product related articles in the Library

Easy chemical structure, query and reaction drawing

Chemical structure drawing

Marvinsketch allows users to quickly draw molecules through basic functions on the GUI and advanced functionalities such as sprout drawing, customizable shortcuts, abbreviated groups, default and user defined templates and context sensitive popup menus.

Atom and Bond properties

Marvinsketch has a rich support for atom and bond properties. Users can assign stereochemistry, charge, valence, radicals and isotopes to each atom. Single, double, triple bonds and aromatic forms are supported. Moreover using wedge bonds user can assign stereochemistry to atoms. Additional data fields can also be attached to atoms, via "S-group" logic so that any user defined information can be stored directly with the structural information.

Reaction drawing

You can draw single step reactions in Marvinsketch by placing a reaction arrow in any position, pointing in any direction in relation to reaction products. The structures 'in front' of the arrow will be recognized as reactants, structures 'above' the arrow as agents, and structures behind as products. Atoms can be automatically or manually mapped using the arrow function.

Query drawing

Atom lists, bond lists, not lists, generic atoms, R-groups, lone pairs are among the query building features available in Marvinsketch. Link nodes, repeating units, pseudo atoms and homology groups, S-groups with attached data are also supported. SMARTS rules allow users to define any specific or generic queries.

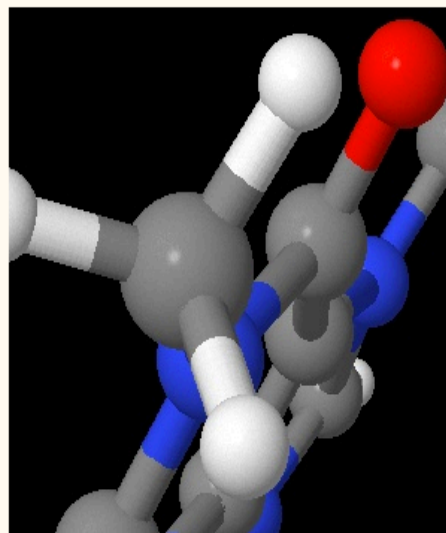
Molecular editors and viewers



<http://jmol.sourceforge.net/>

Jmol: an open-source Java viewer for chemical structures in 3D

with features for chemicals, crystals, materials and biomolecules



Jmol is an interactive web browser applet.

This is a still image, but you can get an animated display of Jmol abilities by clicking [here](#).

(The applet may take some seconds to load. Please, wait and do not reload the page in the meantime.)

The latest stable version is Jmol 12.0

Structure models

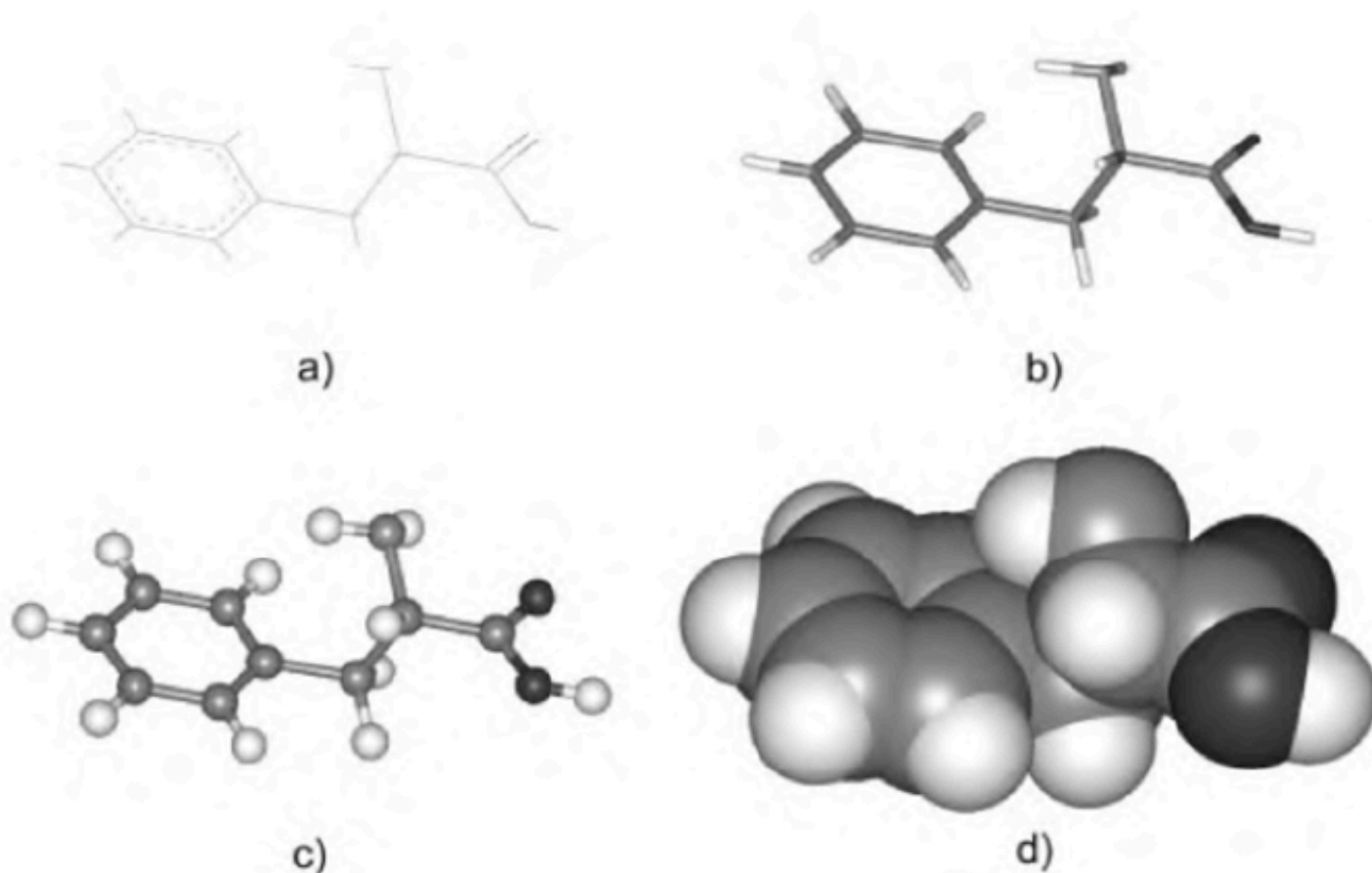


Figure 2-123. The most common molecular graphics representations of phenylalanine a) wire frame; b) capped sticks; c) balls and sticks; d) space-filling.

Format conversion

<http://cactus.nci.nih.gov/translate/>

[Home](#) | [About](#) | [Contact](#) | [Disclaimer](#) | [Privacy](#)



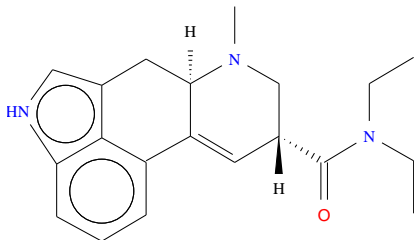
NCI/CADD Group

Online SMILES Translator and Structure File Generator

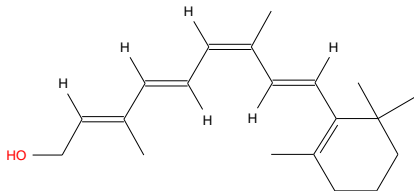
[Form](#) | [News](#) | [Help](#) | [Acknowledgments](#)

Input Format	Unique SMILES Output Format (Unique SMILES)	
<p><input type="text" value="C12C3C4C1C5C4C3C25"/></p> <p><input type="button" value="Start Structure Editor"/></p> <p>Please choose this field if you want to submit your own SMILES strings or create a SMILES string using the Structure Editor. A submitted file has precedence, so delete any entry below if you want to submit a new SMILES string.</p>	<p><input type="radio"/> Display on screen</p> <p><input type="radio"/> SMILES TXT file</p> <p><input type="radio"/> SDF</p> <p><input type="radio"/> PDB</p> <p><input checked="" type="radio"/> MOL (only single structure generated)</p> <p>Use</p> <p><input checked="" type="radio"/> Kekule or</p> <p><input type="radio"/> Aromatic</p> <p>SMILES representation (choose "Aromatic" for closer approximation to Daylight USMILES)</p> <p>SD, PDB or MOL files should contain</p> <p><input checked="" type="radio"/> 2D</p> <p><input type="radio"/> 3D</p> <p>coordinates</p>	
<p><input type="text" value=""/></p> <p><input type="button" value="Browse..."/></p> <p>Please choose this field if you want to translate your own files. The service will automatically recognize SD files (single and multiple structure), text files with multiple SMILES fields, MOL files and PDB files (and in fact any other format CACTVS recognizes).</p>	<p>If the input file contains a single structure, the output will also be single structure. Multiple structure input formats will generate multiple structure output for those formats that support this. Otherwise, only the first structure will be used. SD files will contain a UNIQUE_SMILES field for unique SMILES and an USER_SUPPLIED_SMILES field for the user-supplied SMILES (if available)</p>	
<p><input type="button" value="Reset"/> <input type="button" value="Translate"/></p>		

Examples of SMILES



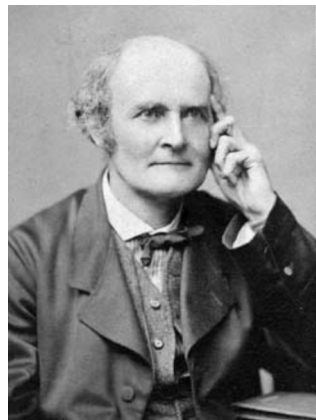
CCN(CC)C(=O)[C@H]1CN(C)[C@@H]2Cc3c[nH]c4cccc(C2=C1)c34



C/C\ (=C\CO)/C=C/C=C(/C)\C=CC1=C(C)CCCC1(C)C

Graph Theory / Algebra / Chemistry

Enumeration of Chemical Isomers



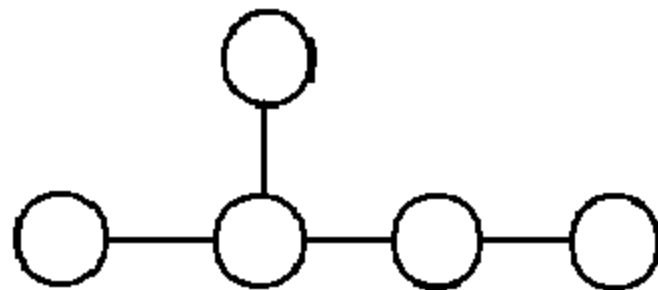
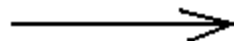
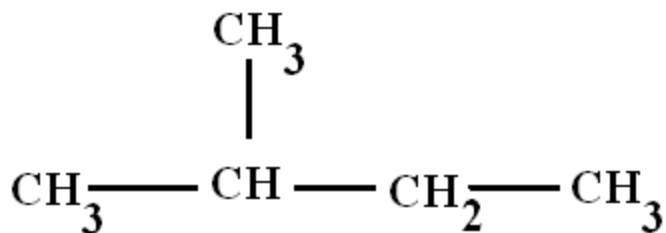
Arthur Cayley



James J. Sylvester

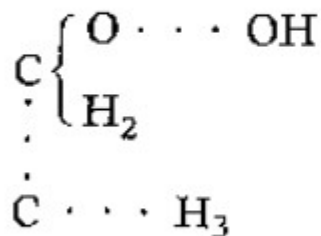


George Polya

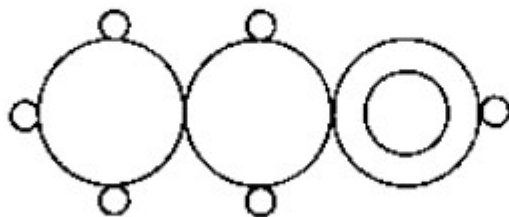


Couper (1858) / Loschmidt (1861) / Kekulé (1861)

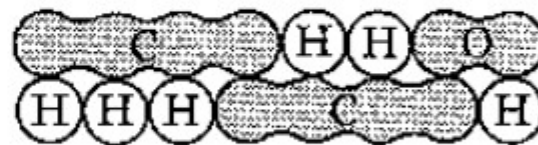
CHEMICAL GRAPHS



Couper

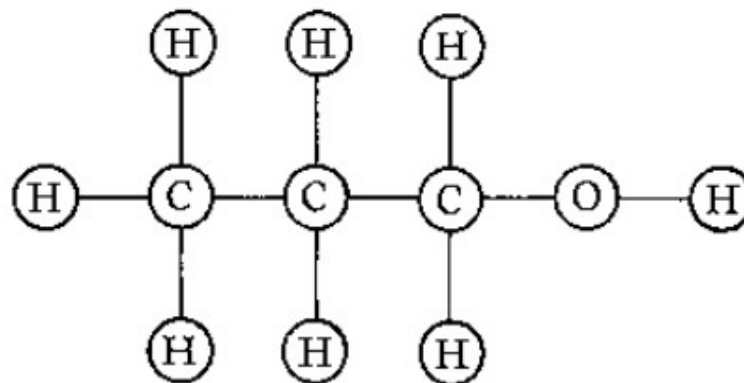
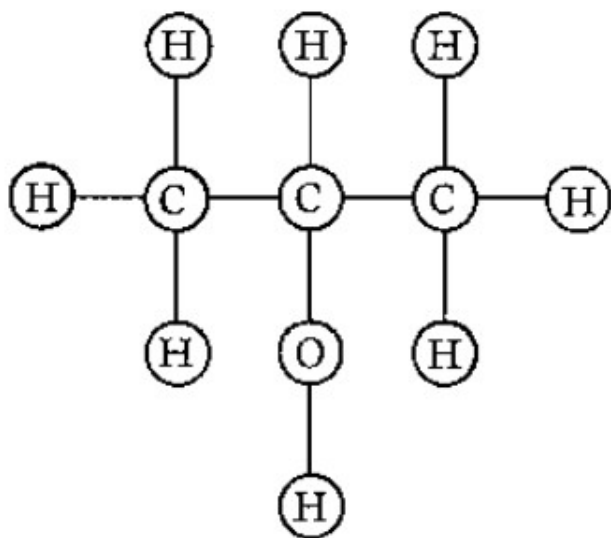


Loschmidt



Kekulé

Crum Brown (1864) and Frankland (1866)



8 cm. from the primary. Reverse the wires in the secondary circuit, reverse the wires in the primary circuit, how you please, the mercury always moves towards the point of the capillary.

8. Shouting or singing (excepting the above-mentioned note) produces no visible effect under the conditions mentioned in Experiments 5, 6, and 7.

9. If the secondary coil be now moved close up, so as to cover as completely as possible the primary, talking to the telephone with the ordinary voice, *i.e.* with moderate strength and at any pitch, produces a definite movement of the mercury column for each word, some sounds of course giving more movement than others, but the movement is always towards the end of the capillary. Singing the note mentioned in Experiments 5, 6, and 7 loudly, produces a movement too large to be measured with the electrometer.

Reversing the poles of the magnet in the telephone does not alter the results of Experiments 5, 6, 7, and 9.

On mentioning the above results to Dr. Burdon Sanderson, he suggested that the apparently anomalous behaviour of the electrometer might be accounted for, by supposing that the mercury moved *quicker* when a current passed towards the point of the capillary than when it flowed in the opposite direction; so that if a succession of rapidly alternating currents be passed through the instrument, the mercury will always move towards the point of the capillary, the movement away from the point being masked by the sluggishness of the instrument in that direction. That this explanation is the correct one is proved by the following experiment:—The current from two Grove's cells is sent through a metal reed vibrating 100 times a second, the contact being made and broken at each vibration, the primary wire of a Du Bois Reymond's induction-coil is also included in the circuit; on connecting the electrometer with the secondary coil placed at an appropriate distance the mercury always moves to the point of the tube whatever be the direction of the current.

F. J. M. PAGE
Physiological Laboratory, University College,
London, February 2

NOTE.—On February 4 Prof. Graham Bell kindly placed at my disposal a telephone much more powerful than any of those I had previously used. On speaking to this instrument, the electrometer being in the circuit, movements of the mercury column as considerable as those in Experiment 9 were observed.—F. J. M. P.

CHEMISTRY AND ALGEBRA

IT may not be wholly without interest to some of the readers of NATURE to be made acquainted with an analogy that has recently forcibly impressed me between branches of human knowledge apparently so dissimilar as modern chemistry and modern algebra. I have found it of great utility in explaining to non-mathematicians the nature of the investigations which algebraists are at present busily at work upon to make out the so-called *Grundformen* or irreducible forms appurtenant to binary quantics taken singly or in systems, and I have also found that it may be used as an instrument of investigation in purely algebraical inquiries. So much is this the case that I hardly ever take up Dr. Frankland's exceedingly valuable "Notes for Chemical Students," which are drawn up exclusively on the basis of Kekulé's exquisite conception of *valence*, without deriving suggestions for new researches in the theory of algebraical forms. I will confine myself to a statement of the grounds of the analogy, referring those who may feel an interest in the subject and are desirous for further information about it to a memoir which I have written upon it for the new *American Journal of Pure and Applied Mathematics*, the first number of which will appear early in February.

The analogy is between atoms and binary quantics exclusively.

I compare every binary quantic with a chemical atom. The number of factors (or rays, as they may be regarded by an obvious geometrical interpretation) in a binary quantic is the analogue of the number of *bonds*, or the *valence*, as it is termed, of a chemical atom.

Thus a linear form may be regarded as a monad atom, a quadratic form as a duad, a cubic form as a triad, and so on.

An invariant of a system of binary quantics of various degrees is the analogue of a chemical substance composed of atoms of corresponding *valences*. The order of such invariant in each set of coefficients is the same as the number of atoms of the corresponding *valence* in the chemical compound.

A co-variant is the analogue of an (organic or inorganic) compound radical. The orders in the several sets of coefficients corresponding, as for invariants, to the respective valences of the atoms, the free valence of the compound radical then becomes identical with the degree of the co-variant in the variables.

The weight of an invariant is identical with the number of the bonds in the chemicograph of the analogous chemical substance, and the weight of the leading term (or basic differentiant) of a co-variant is the same as the number of bonds in the chemicograph of the analogous compound radical. Every invariant and covariant thus becomes expressible by a *graph* precisely identical with a Kekuléan diagram or chemicograph. But not every chemicograph is an algebraical one. I show that by an application of the algebraical law of reciprocity every algebraical graph of a given invariant will represent the constitution in terms of the roots of a quantic of a type reciprocal to that of the given invariant of an invariant belonging to that reciprocal type. I give a rule for the geometrical multiplication of graphs, *i.e.* for constructing a *graph* to the product of in- or co-variants whose separate graphs are given. I have also ventured upon a hypothesis which, whilst in nowise interfering with existing chemicographical constructions, accounts for the seeming anomaly of the isolated existence as "monad molecules" of mercury, zinc, and arsenic—and gives a rational explanation of the "mutual saturation of bonds."

I have thus been led to see more clearly than ever I did before the existence of a common ground to the new mechanism, the new chemistry, and the new algebra. Underlying all these is the theory of pure colligation, which applies undistinguishably to the three great theories, all initiated within the last third of a century or thereabouts by Eisenstein, Kekulé, and Peaucellier.

Baltimore, January 1

J. J. SYLVESTER

PALMEN ON THE MORPHOLOGY OF THE TRACHEAL SYSTEM

DR. PALMEN, of Helsingfors, has recently published an interesting memoir on the tracheal system of insects. He observes that although the gills of certain aquatic larvæ are attached to the skin very near to the points at which the spiracles open in the mature insects, and though spiracles and gills do not co-exist in the same segment, yet the point of attachment of the gills never exactly coincides with the position of the future spiracle. Moreover, he shows that even during the larval condition, although the spiracles are not open, the structure of the stigmatic duct is present, and indeed that it opens temporarily at each moult, to permit the inner tracheal membrane to be cast, after which it closes again. In fact, then, he urges, the gills and spiracles do not correspond exactly, either in number or in position, and there can therefore be between them no genetic connection. He concludes that the insects with open tracheæ are not derived from ancestors provided with gills,

Some Polya Enumeration

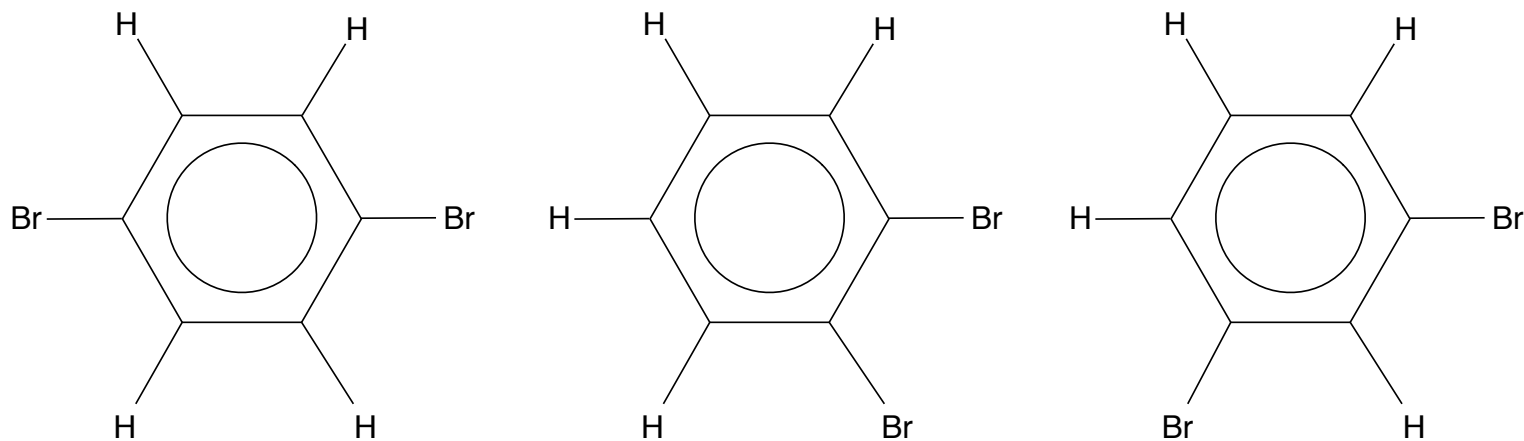
Pólya Theory

Formally: Pólya Theory counts equivalence classes, where the equivalence classes are induced by group actions. Since groups describe symmetry, Pólya Theory is counting the number of distinct objects in the presence of symmetry.

Informally: Pólya Theory does “common sense” counting.

Chemical Isomers

Pólya's original objective was to determine the number of distinct compounds given a chemical formula.



There are three distinct compounds with the formula $C_6H_4Br_2$.

Example

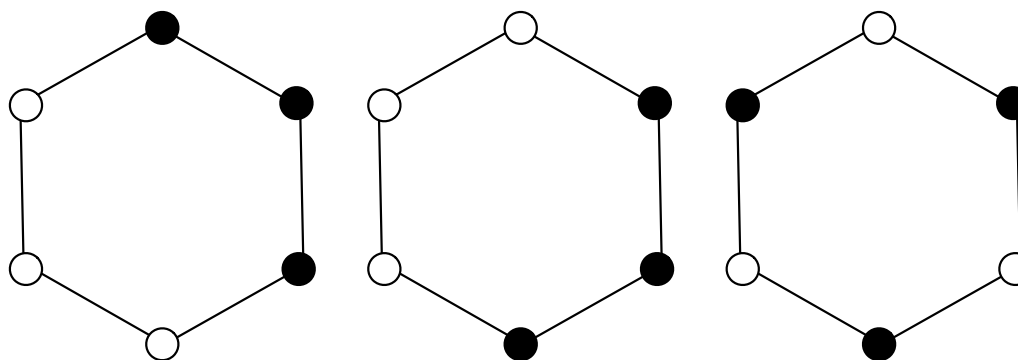


Figure 1: Common sense says that these the two bracelets on the left are the “same”, the third bracelet is “different”

Another Example

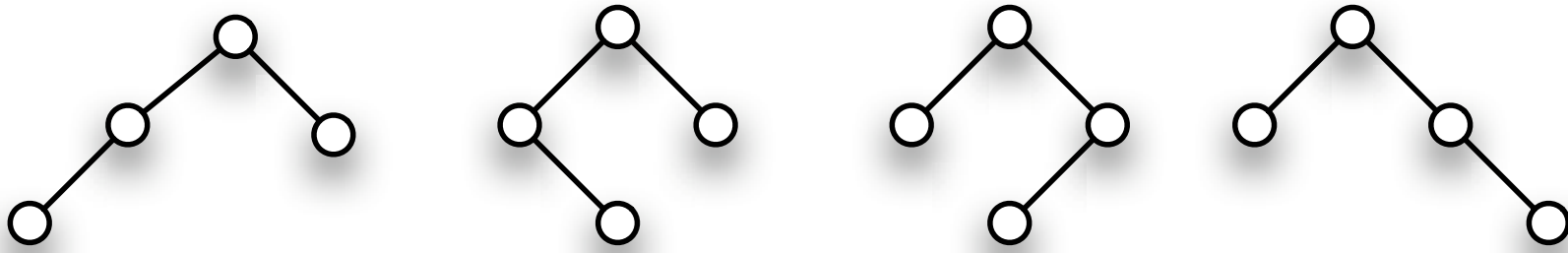


Figure 2: These 4 binary trees are equivalent if left and right are considered indistinguishable

Yet Another Example

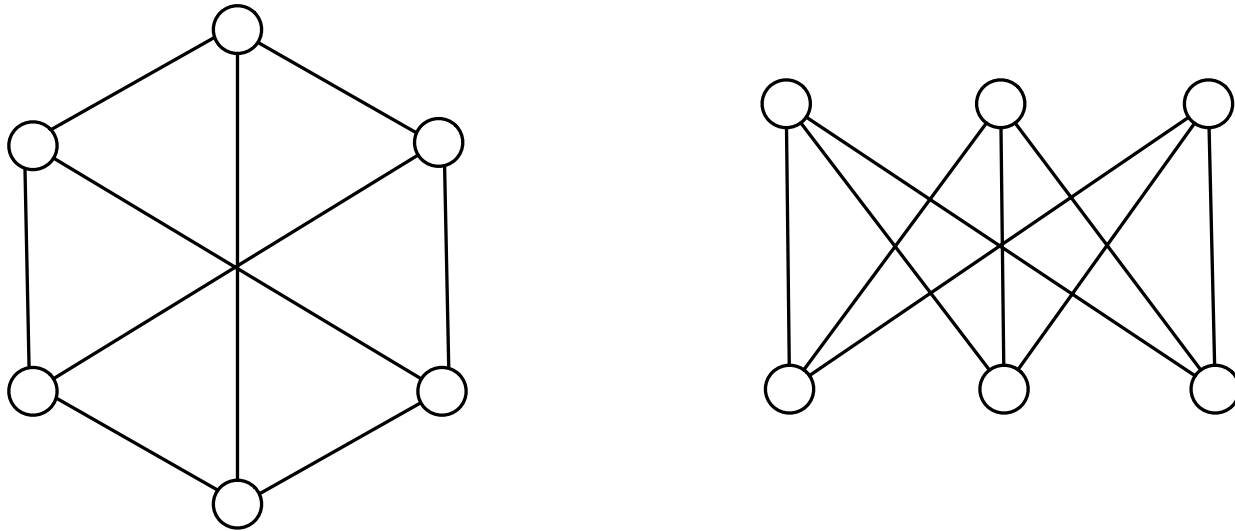


Figure 3: These two graphs are isomorphic, although it is not visually obvious

And Another

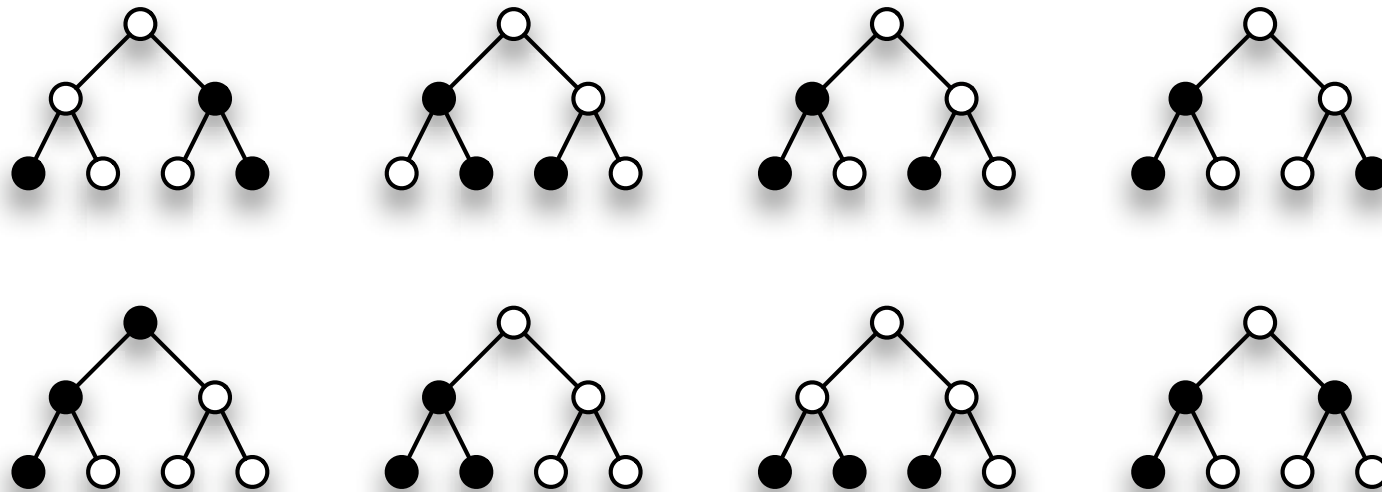
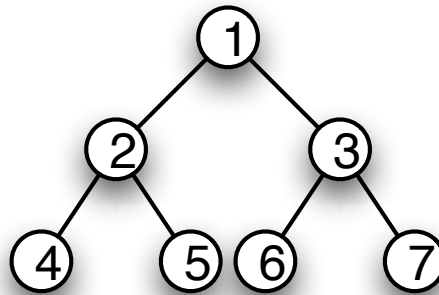


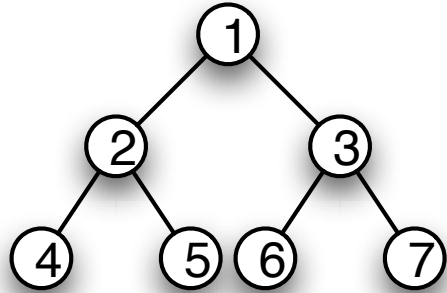
Figure 4: If left and right are indistinguishable then the top row of graph colorings are equivalent, and the bottom row are different.

Example Calculation

Problem: Count the number of distinct black and white colorings of the tree in Figure 7 if left and right are indistinguishable.



Step 1



Write down the permutations that leave the tree invariant if left and right are indistinguishable. This is a group G generated by:

$$\pi_1 = (1)(2\ 3)(4\ 6)(5\ 7)$$

$$\pi_2 = (1)(2)(3)(4\ 5)(6)(7)$$

$$\pi_3 = (1)(2)(3)(4)(5)(6\ 7)$$

Step 2

List the group elements, calculate the monomials.

Element	Cycle Representation	Monomial
I	$(1)(2)(3)(4)(5)(6)(7)$	x_1^7
π_1	$(1)(2\ 3)(4\ 6)(5\ 7)$	$x_1 x_2^3$
π_2	$(1)(2)(3)(4\ 5)(6)(7)$	$x_1^5 x_2$
π_3	$(1)(2)(3)(4)(5)(6\ 7)$	$x_1^5 x_2$
$\pi_2 \pi_3 = \pi_3 \pi_2$	$(1)(2)(3)(4\ 5)(6\ 7)$	$x_1^3 x_2^2$
$\pi_1 \pi_2 = \pi_3 \pi_1$	$(1)(2\ 3)(4\ 6\ 5\ 7)$	$x_1 x_2 x_4$
$\pi_1 \pi_3 = \pi_2 \pi_1$	$(1)(2\ 3)(4\ 7\ 5\ 6)$	$x_1 x_2 x_4$
$\pi_1 \pi_2 \pi_3$	$(1)(2\ 3)(4\ 7)(5\ 6)$	$x_1 x_2^3$

Table 1: The elements of G

Step 3

Add the monomials to get the *cycle index* for the group:

$$P_G(x_1, x_2, x_4) = \frac{1}{8}(x_1^7 + 2x_1^5x_2 + 2x_1x_2^3 + 2x_1x_2x_4 + x_1^3x_2^2)$$

Formally the cycle index is defined:

$$P_G(x_1, x_2, \dots, x_{|D|}) = \frac{1}{|G|} \sum_{\pi \in G} x_1^{l_1(\pi)} x_2^{l_2(\pi)} \dots x_{|D|}^{l_{|D|}(\pi)}$$

where D is the the set acted on by elements of G , $|D|$ is the size of the set, $l_k(\pi)$ is the number of cycles of length k in π .

Step 4

Pólya's Enumeration Theorem says that the number of distinct k -colorings is

$$\begin{aligned} P_G(k, k, k, k) &= \frac{1}{8}(k^7 + 2k^6 + 2k^4 + 2k^3 + k^5) \\ &= \frac{k^3}{8}(k^4 + 2k^3 + k^2 + 2k + 2) \\ &= \frac{k^3}{8}(k + 1)(k^3 + k^2 + 2) \end{aligned}$$

This must be an integer for all integer values of k , so $k^3(k + 1)(k^3 + k^2 + 2)$ must be divisible by 8.

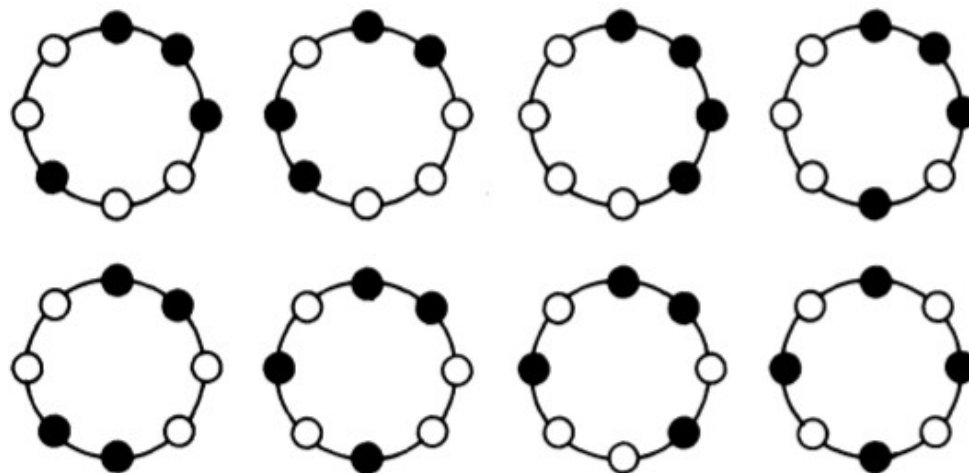
The Solution

The number of 2-colorings of the binary tree, with left and right indistinguishable, is

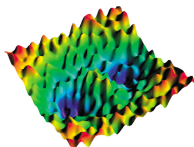
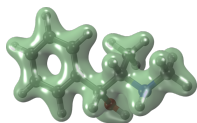
$$\begin{aligned} P_G(2, 2, 2, 2) &= \frac{1}{8}(2^7 + 2^5 + 2^7 + 2^4 + 2^5) \\ &= 42 \end{aligned}$$

$$1 + x + 4x^2 + 5x^3 + 8x^4 + 5x^5 + 4x^6 + x^7 + x^8.$$

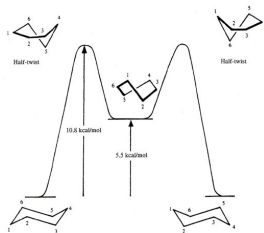
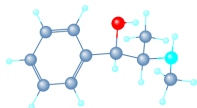
Thus, for example, there are five inequivalent necklaces having three black beads, and eight with equal numbers of black and white beads. The latter are shown in FIGURE 1.



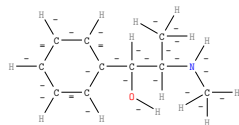
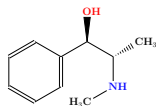
Levels of Abstraction in Computational Chemistry



Potential energy surface



Reaction coordinate



Graph grammar

[Andersen et al., Proceedings of the Royal Society A, 2017]

Levels of Abstraction in Programming

The screenshot shows the Compiler Explorer interface. On the left, the C++ source code for a function named `testFunction` is displayed. The code calculates the sum of an array of integers. On the right, the corresponding assembly code for x86-64 gcc 6.3 is shown. The assembly code includes instructions for stack frame setup, loop initialization, and the loop body logic, demonstrating the low-level implementation of the high-level C++ code.

```
1 int testFunction(int* input, int length) {
2     int sum = 0;
3     for (int i = 0; i < length; ++i) {
4         sum += input[i];
5     }
6     return sum;
7 }
8
```

```
11010 .L0: .text // Intel
1 testFunction(int*, int):
2     push    rbp
3     mov     rbp, rsp
4     mov     QWORD PTR [rbp-24], rdi
5     mov     DWORD PTR [rbp-28], esi
6     mov     DWORD PTR [rbp-4], 0
7     mov     DWORD PTR [rbp-8], 0
8     .L3:
9     mov     eax, DWORD PTR [rbp-8]
10    cmp     eax, DWORD PTR [rbp-28]
11    jge    .L2
12    mov     eax, DWORD PTR [rbp-8]
13    cdqeq
14    lea    rdx, [0rax*4]
15    mov     rax, QWORD PTR [rbp-24]
16    add     rax, rdx
17    mov     eax, DWORD PTR [rax]
18    add     DWORD PTR [rbp-4], eax
19    add     DWORD PTR [rbp-8], 1
20    jmp    .L3
21    .L2:
22    mov     eax, DWORD PTR [rbp-4]
23    pop     rbp
24    ret
```

Declarative Description ↔ DSL ↔ C++ ↔ Assembler

Levels of Abstraction in Computer Science



*“The psychological profiling [of a Computer Scientist] is mostly the ability to **shift levels of abstraction**, from low level to high level. To see something in the small and to see something in the large.”*

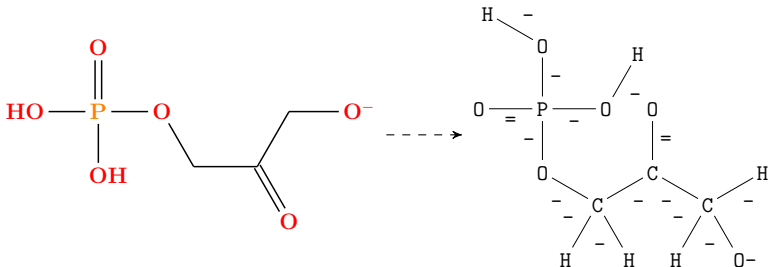
Donald Knuth

Modelling and Analysis of Chemical Systems

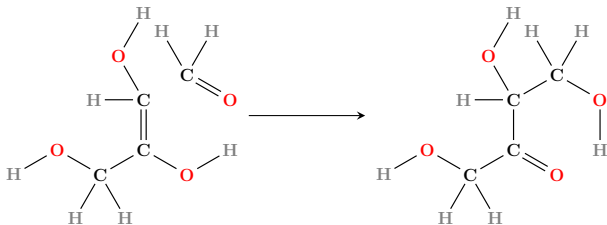
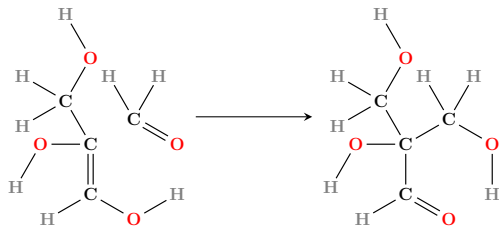
Modelling and Analysis of Chemical Systems

1. Model molecules as labelled graphs.

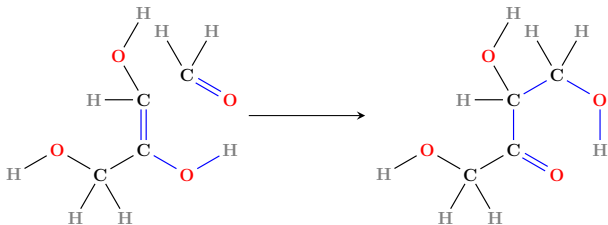
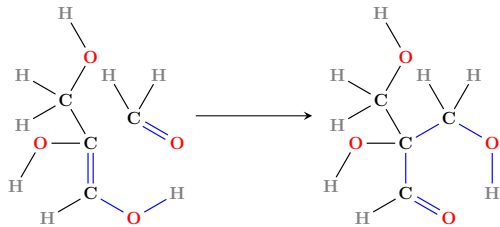
- ▶ **An old idea:** [J. J. Sylvester, *Chemistry and Algebra*, Nature 1878]
- ▶ **Molecule:** simple, connected, labelled graph.
- ▶ **Vertex labels:** atom type, charge.
- ▶ **Edge labels:** bond type.



Chemical Reactions (Educts \rightarrow Products)

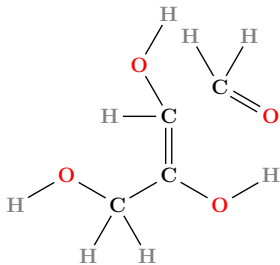
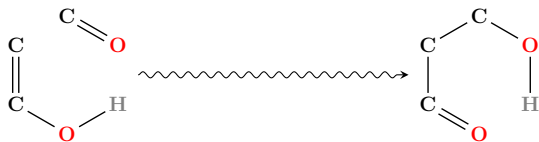


Chemical Reactions (of the Same Type)



Chemical Reaction Patterns

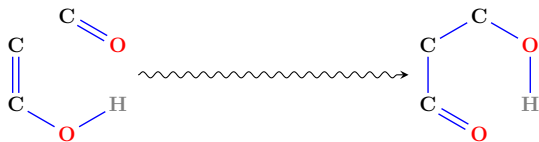
Rule



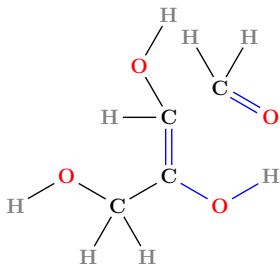
Educts

Chemical Reaction Patterns

Rule

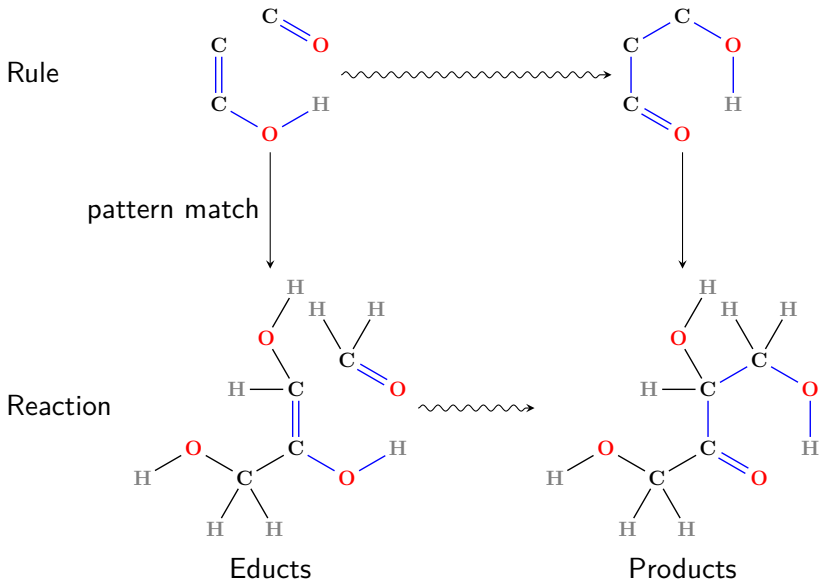


pattern match



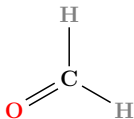
Educts

Chemical Reaction Patterns

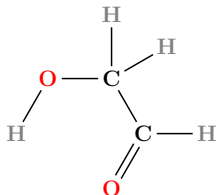


Grammar Example: The Formose Chemistry

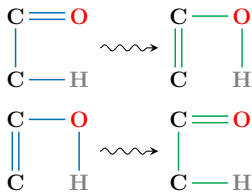
Formaldehyde:



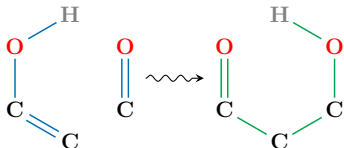
Glycolaldehyde:



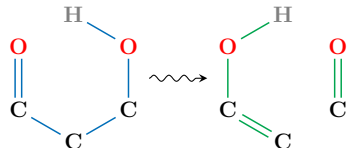
Keto-enol tautomerism:



Aldol addition:



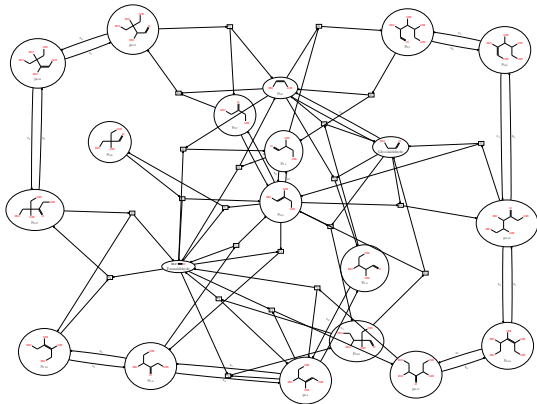
Retro aldol addition:



Modelling and Analysis of Chemical Systems

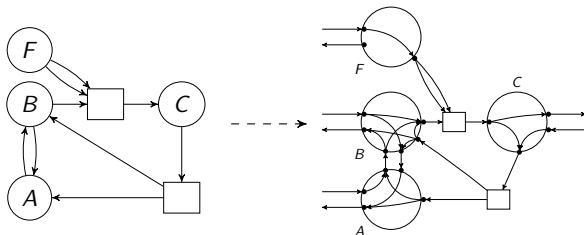
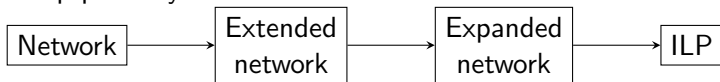
3. Generate a reaction network.

```
dg = dgRuleComp(inputGraphs ,  
  addSubset(inputGraphs) >> rightPredicate[  
    lambda d: all(countCarbon(a) <= 5 for a in d.right)  
  ](  
    repeat(inputRules)    )  
  )  
dg.calc()
```



Modelling and Analysis of Chemical Systems

4. Set up pathway model.



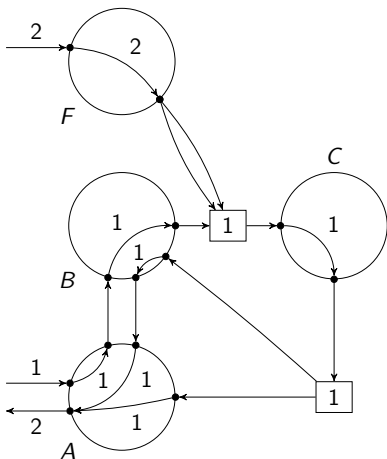
Conservation constraints:

$$\sum_{e \in \delta_E^+(v)} m_v(e^+) f(e) - \sum_{e \in \delta_E^-(v)} m_v(e^-) f(e) = 0 \quad \forall v \in \tilde{V}$$

Modelling and Analysis of Chemical Systems

5. Formulate pathway question.

Example: Given 2 formaldehyde and 1 glycolaldehyde, how can 2 glycolaldehyde be produced through autocatalysis.



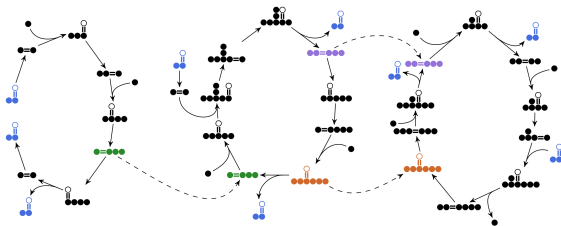
Modelling and Analysis of Chemical Systems

6. Enumerate many alternate pathways.

Example (Formose):

Network: all molecules with at most 9 carbon atoms.

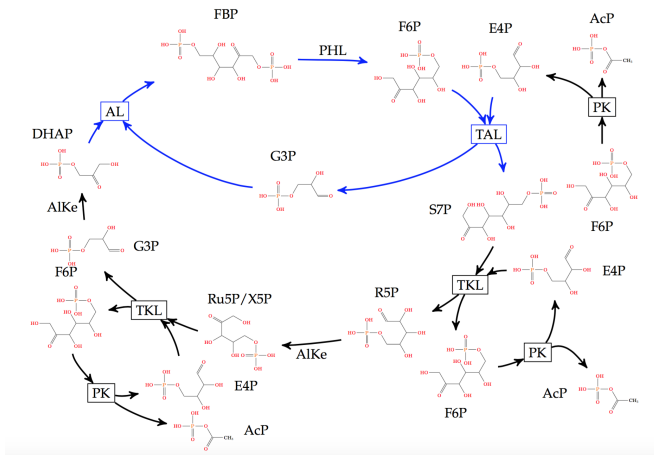
Reactions used	Maximum #C						Sum
	4	5	6	7	8	9	
6	0	0	1	1	1	2	5
7	0	0	0	0	0	2	2
8	1	5	7	17	37	68	135
9	0	0	12	12	37	69	130
10	0	12	50	274	849	—	≥ 1185
11	0	5	41	190	738	—	≥ 974
							≥ 2431



Another Example: Non-oxidative Glycolysis

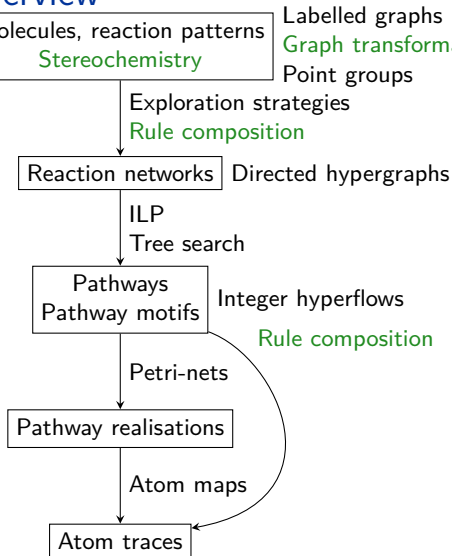
You specify: $F6P + 2 P_i \rightarrow 3 AcP + 2 H_2O$

You get (for example):



Many alternatives for suggestion in [Bogorad, Lin, and Liao, Nature, 2013]

Overview



Labelled graphs
Graph transformation rules
Point groups

Category theory
Double Pushout
Rule composition
Monomorphisms
Isomorphisms
Canonicalisation
Automorphisms

Software package: MØD
C++, Python, Bash, L^AT_EX

Pentose phosphate pathway
Glycolysis (EMP and ED)
Non-oxidative glycolysis
Citric acid cycle
Enzyme mechanisms
Formose
Prebiotic chemistry (HCN)
Eschenmoser's GLX scenario
DNA templated computing

<http://cheminf.imada.sdu.dk>