

Graph algorithms in crystallography and cheminformatics

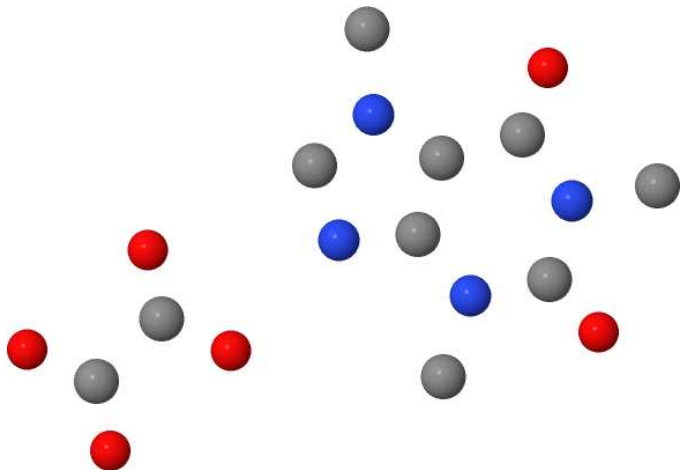
Andrius Merkys

Grenoble, 2024

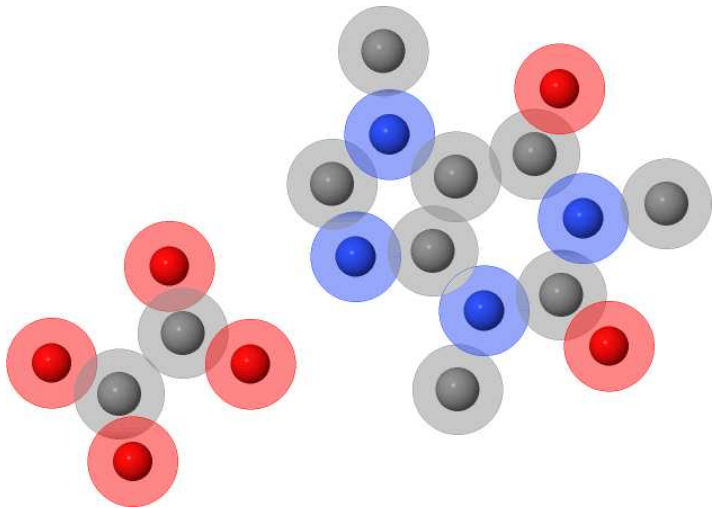
Vilnius University Institute of Biotechnology



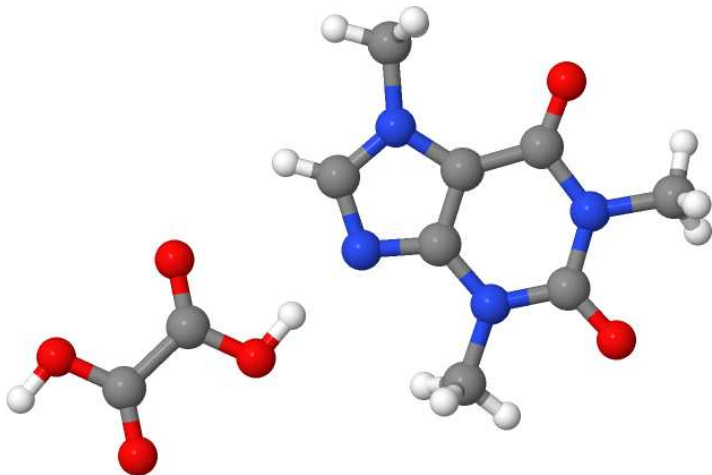
Output of X-ray crystallography



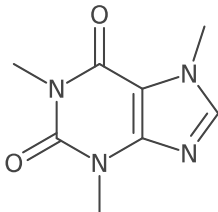
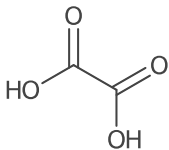
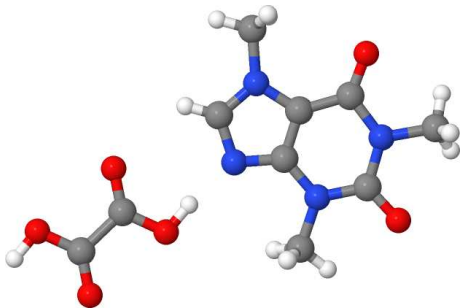
Detection of connectivity



Molecular graph + 3D coordinates



Structural formula



Properties of molecular graphs

Properties

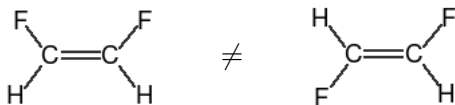
- ▶ undirected
- ▶ no self-loops
- ▶ non-multigraphs (no parallel edges)
- ▶ usually planar (Simmons & Maggio, 1981)

Colors

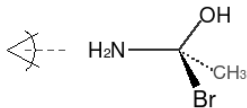
- ▶ vertices – chemical element, sometimes: charge, isotope
- ▶ edges – bond order (1, 2, 3, 4)

Stereochemistry

- ▶ Cis/trans configuration around double bonds



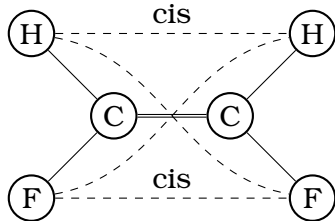
- ▶ Tetrahedral chirality



<http://opensmiles.org/opensmiles.html>
OpenSMILES specification

Cis/trans configuration in molecular graphs

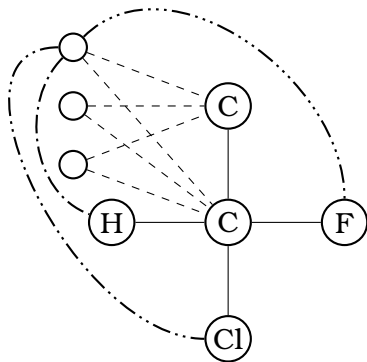
- ▶ Spatial relations can be represented by additional colored edges



Merkys et al., 2023

Tetrahedral chirality in molecular graphs

- ▶ Some relations need additional vertices



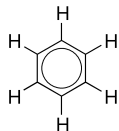
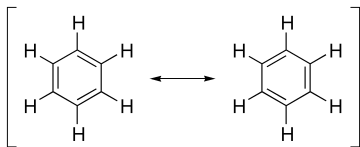
Merkys et al., 2023

Cycles

Importance

- ▶ hold the structure (limit the conformational space)
- ▶ base for aromaticity

Aromatic cycles



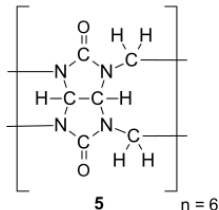
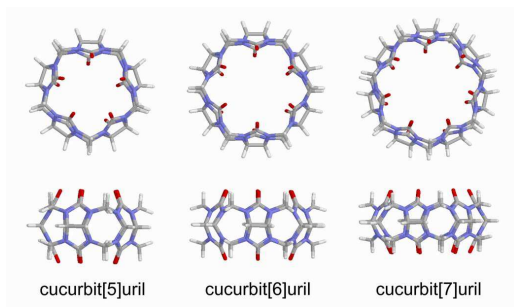
Smallest set of smallest rings (SSSR)

Closed path \in SSSR iff there are no chords

Downs et al. 1989 algorithm:

1. choose L – length of the largest cycle
2. find all paths from all vertices (depth-first search)
3. search is terminated when L is reached
4. search is terminated when other criteria for SSSR are violated

Cycles. Example: cucurbituril



- ▶ $n \times 2$ five-membered cycles
- ▶ n eight-membered cycles
- ▶ two $n \times 4$ -membered cycles

https://commons.wikimedia.org/wiki/File:Models_of_cucurbiturils.jpg

M stone at English Wikipedia, CC-BY-SA 3.0

<https://commons.wikimedia.org/wiki/File:CucurbiturilSynthesis.svg>

V8rik at English Wikipedia, CC-BY-SA 3.0

All cycles in a graph

All closed paths without same vertex appearing more than once

Hanser, Jauffret, Kaufmann 1996 algorithm:

1. choose vertex v and remove it from graph
2. connect v neighbours with edges representing paths (through v) between them
3. every self-loop found is a cycle in the original graph
4. loop to step 1 until graph becomes empty

All cycles in a graph (2)

Examples

- ▶ cube – 28 cycles
- ▶ fullerene (a.k.a. C_{60}) – 8018 cycles

Problems

- ▶ computationally expensive
- ▶ no known low-level (C/C++) implementations

Search and comparison

Search – look for a matching subgraph

- ▶ boils down to subgraph isomorphism problem
- ▶ property vectors (a.k.a. fingerprints) may be used

Comparison – check if two graphs match

- ▶ identity – graph isomorphism problem
- ▶ property vectors (a.k.a. fingerprints) may be used

Property vectors a.k.a. fingerprints

Tanimoto similarity index

$$T_s(A, B) = \frac{\sum_i A_i \wedge B_i}{\sum_i A_i \vee B_i}$$

Example:

$$T_s(1001, 1010) = \frac{1}{3}$$

Open Babel v3.1.1 uses:

- ▶ FP3 – 55 properties
- ▶ FP4 – 309 properties

Molecular graph representation formats

- ▶ IUPAC preferred name (*Pyridine*)
 - ▶ uniquely defines a compound
 - ▶ difficult to read/write by computer
- ▶ SMILES (c1ccncc1)
 - ▶ easy to read/write by computer
 - ▶ could be read/write by human
 - ▶ many competing standards
 - ▶ not a unique representation
- ▶ InChI (InChI=1S/C5H5N/c1-2-4-6-5-3-1/h1-5H)
 - ▶ uniquely defines a compound
 - ▶ difficult to read/write by human
 - ▶ sole I/O library is non-free software*

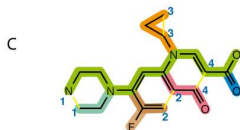
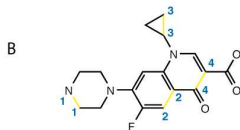
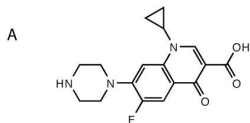
IUPAC names

- ▶ Reading
 - ▶ **OPSIN**
- ▶ Writing
 - ▶ ChemAxon
 - ▶ Lexichem
 - ▶ Nomenclator
 - ▶ STOUT
 - ▶ ChemOnomatopist

	ChemOnomatopist v0.9.0	STOUT v2.0
Correct PIN	1262	1132
Alternative name	752	1874
Incorrect	1074	690
Refused	608	0
Time	≈ 4 min.	≈ 230 min.

SMILES format

- ▶ depicts spanning tree of a graph
- ▶ vertices are named after chemical elements
- ▶ branches from main chain are written in parentheses
- ▶ edges not in spanning tree are marked by numbers
- ▶ bonds of order 2 or more are depicted by =, #, \$



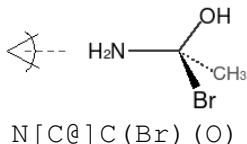
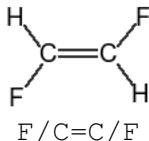
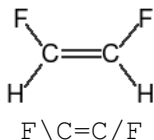
D

N1CCN(CC1)C(C(F)=C2)=CC(=C2C4=O)N(C3CC3)C=C4C(=O)O

<https://commons.wikimedia.org/wiki/File:SMILES.png>
Fdardel & DMacks, CC-BY-SA 3.0

SMILES – representation of a molecular graph

- ▶ aromatic atoms are written in lowercase
- ▶ cis/trans configuration is marked with up/down edges
- ▶ tetrahedral chirality is marked with @ or @@ and establishing atom enumeration order



- ▶ support for chirality for higher coordination numbers

<http://opensmiles.org/opensmiles.html>
OpenSMILES specification

Morgan's algorithm

Steps

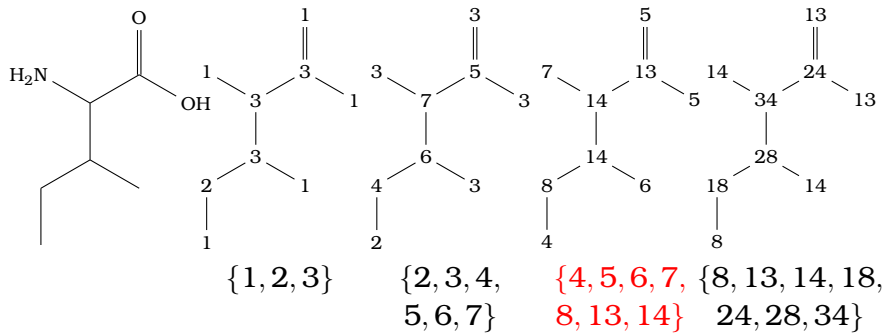
1. $S_0(v) = \text{deg}(v), \forall v \in V$
2. $S_i(v) = \sum S_{i-1}(n)$, sum over all neighbours of v
3. ...

Loop is stopped when the number of distinct $S_i(v)$ stops increasing.

Equivalence classes are defined by $S_i(v)$ values.

Morgan, 1965

Morgan's algorithm. Example

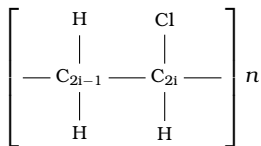


Molecular graph isomorphism

- ▶ Faulon's algorithm (Faulon, 1998)
 - ▶ input – simple graph
 - ▶ $O(N^2)$, observed on carbon nanotubes

- ▶ *Nauty* (McKay & Piperno, 2014)
 - ▶ supports vertices with attributes (color)
 - ▶ employed in InChI
 - ▶ $O(N)$ (Faulon, Collins & Carr, 2004)

Polymeric molecules



Construction of molecular graph

1. choose a representing “monomer” (how?)
2. redirect boundary-crossing edges back to the “monomer”
 - ▶ parallel edges or self-loops may appear
 - ▶ multigraphs may be represented by line graphs

Representing polymer molecules with quotient graphs

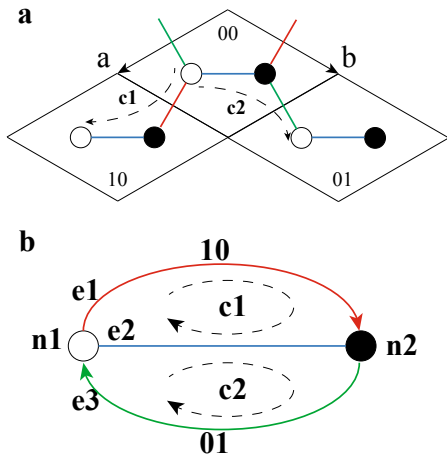
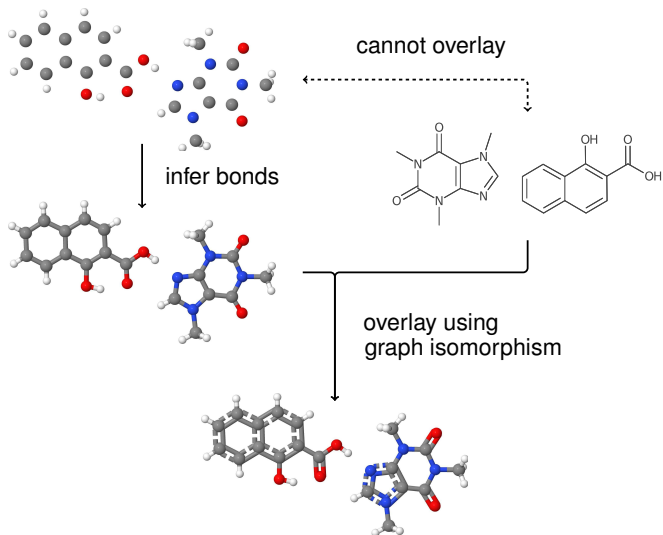


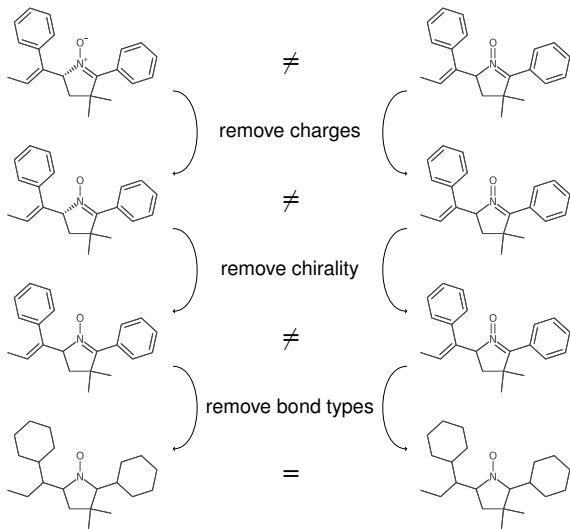
Fig. 1 Graphene net and its quotient graph. a Graphene net in 2D space; b The quotient graph of graphene net. n_1 , n_2 are nodes. e_1 , e_2 , e_3 are edges. c_1 , c_2 are the two basic cycles of graphene net.

Chemical annotation of X-ray structures



<https://www.crystallography.net/archives/2021/posters/IUCr-XXV/poster-1476.pdf>
Merkys et al., poster presentation at IUCr25

Comparison of molecular graphs



<https://www.crystallography.net/archives/2021/posters/IUCr-XXV/poster-1476.pdf>
Merkys et al., poster presentation at IUCr25

Author-provided vs. curated annotations in COD

#	H atoms	aromaticity	atom types	charge	chirality	cis/trans	order	extra moieties
18424		x					x	
2990		x				x	x	
2803		x		x			x	
2191								
1878		x			x		x	
754					x			
467	x	x		x			x	
464		x		x		x	x	
381	x			x				
270		x					x	x
250	x	x					x	
192		x		x	x		x	
164		x			x	x	x	
145						x		
116				x			x	
85	x		x				x	
78								x
75	x	x				x	x	
56			x				x	
45				x				
37	x						x	
36	x							
35		x		x			x	x
34	x			x			x	
33	x	x		x			x	x
32474	<i>total (471 not shown here for brevity)</i>							
2196	unknown							

Problems

- ▶ What is a bond?
- ▶ How to describe interactions between more than two atoms?
- ▶ How to uniquely represent aromatic cycles?