

A chemical notation system: SMILES

Goals:

- (1) uniqueness of the description of the molecule graph
- (2) user-friendly / human readability
- (3) machine-friendly

Rules:

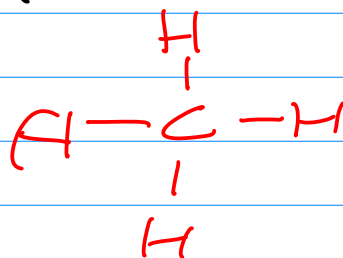
SMILES notation: a series of characters

- (1) atoms represented by their atomic symbol, enclosed in square brackets

elements of the organic subset may be written without explicit hydrogens

C, N, O, P, B, S, F, Cl, Br

C



(2) bonds

single bond

—

double bond

=

triple bond

#

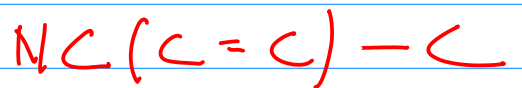
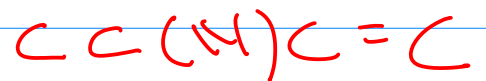
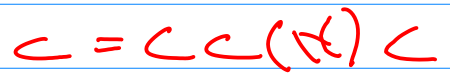
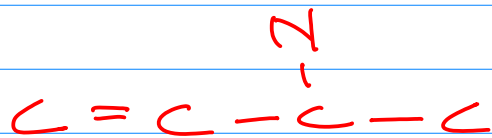
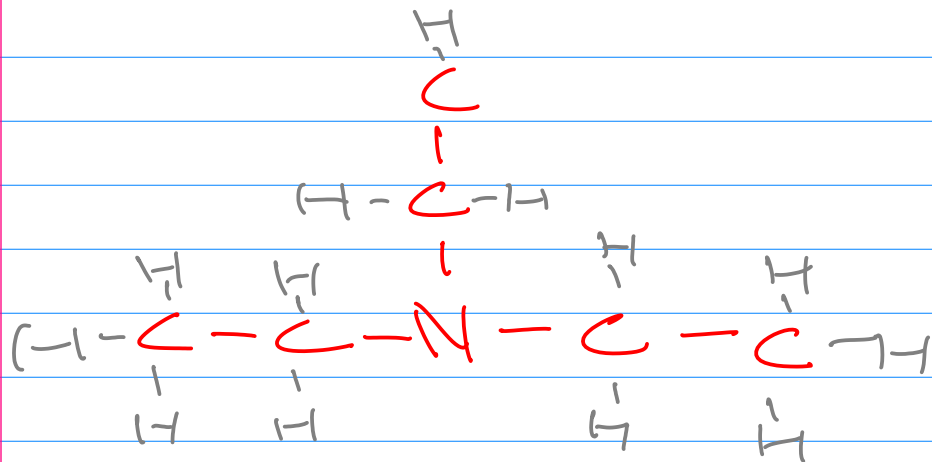
aromatic bond

:

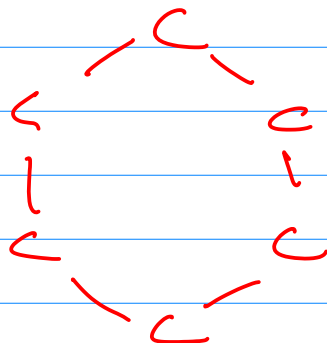
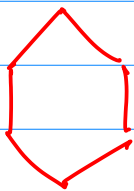
may be omitted

③) Structures

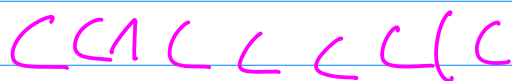
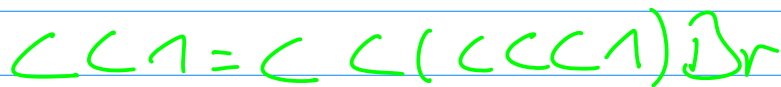
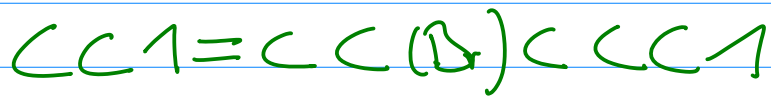
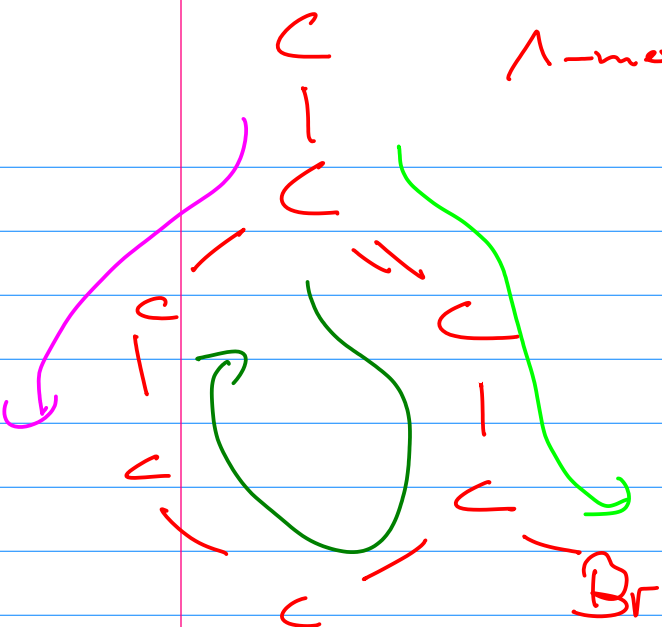
encoded by enclosure parentheses



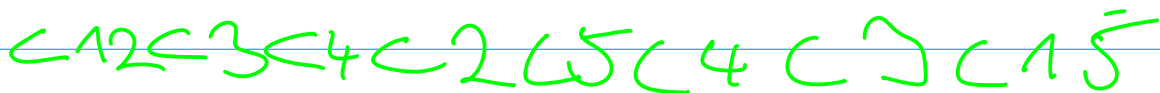
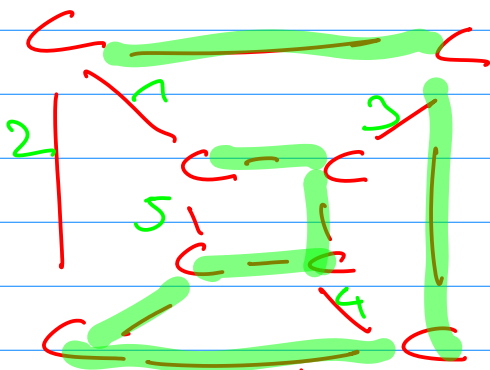
(4) cyclic structures



1-methyl-3-bromo-cyclohexene

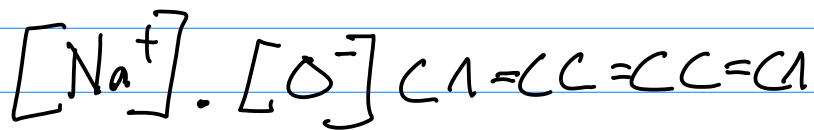
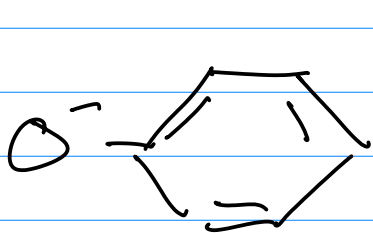


single atom might have many dashes:
cuds are

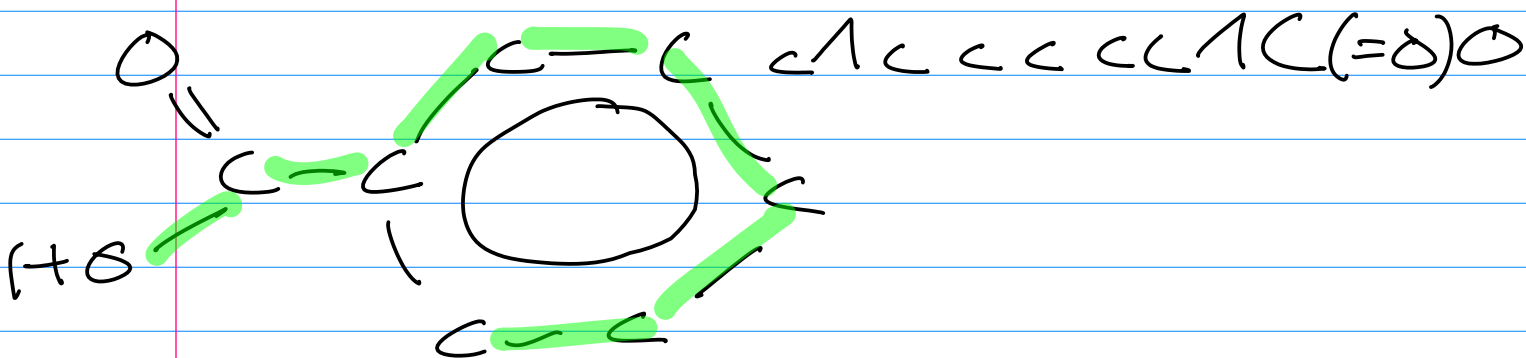


5) discretized structures

Na⁺



G) aromatic rings → small letters
benzoic acid



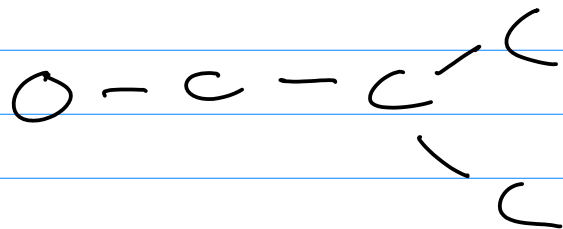
Generation of Unique SMILES Notation

method **CANON** and **GENES**

canonical labeling

generation of unique SMILES

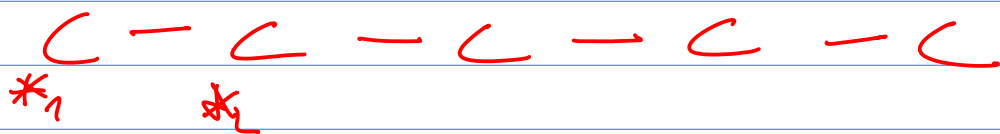
CANON must identify topologically equivalent in all respects molecules



CANON

(a) initial graph is canonized

- (1) # connections
- (2) # non-H bonds
- (3) atomic number
- (4) sign of charge
- (5) absolute charge
- (6) # attached hydrogens



*₁ (methyl carbon): 1, 1, 6, 0, 0, 3

*₂ (methylene carbon): 2, 2, 6, 0, 0, 3

116003 - 226003 - 226003 - ...

(b) rank equivalence

1-2-2-2-1

(c) simple extended connectivity

replace rank by (rank, # rank neighbors)

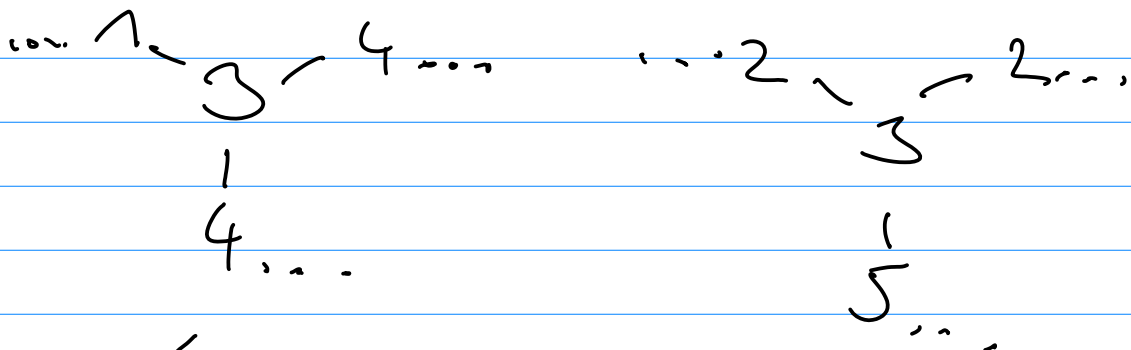
(1,2) - (2,3) - (2,4) - (2,3) - (1,2)

and replace:

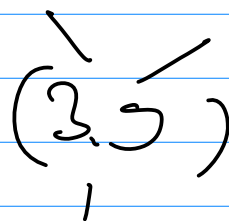
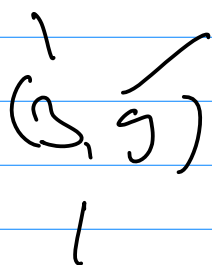
1-2-3-2-1

(d) extended connectivity using unambiguous functions

Why:



↳

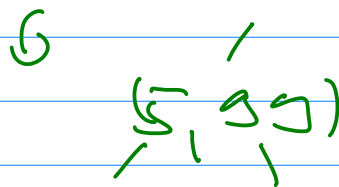
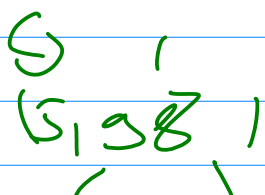
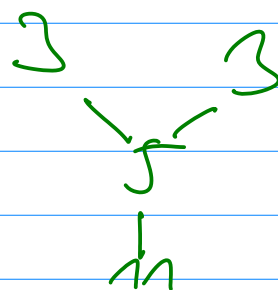
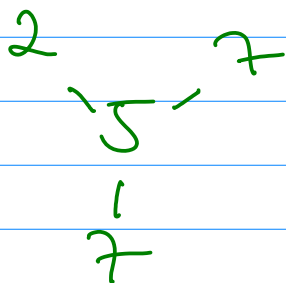


Solution: - replace rules by primes

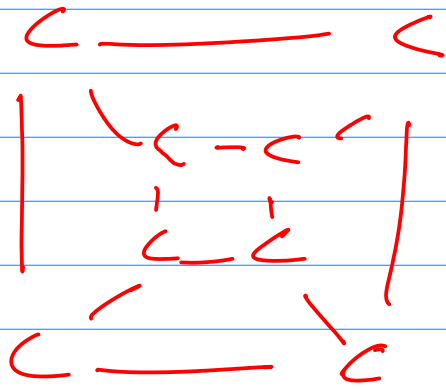
2, 3, 5, 7, 11, ...

- use product instead of sum

↳



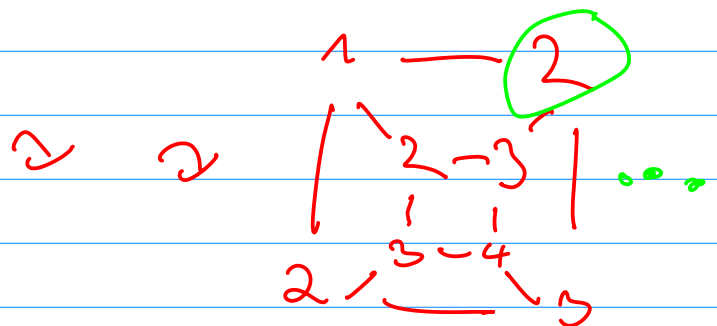
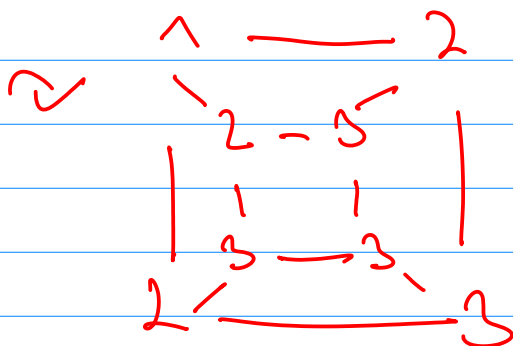
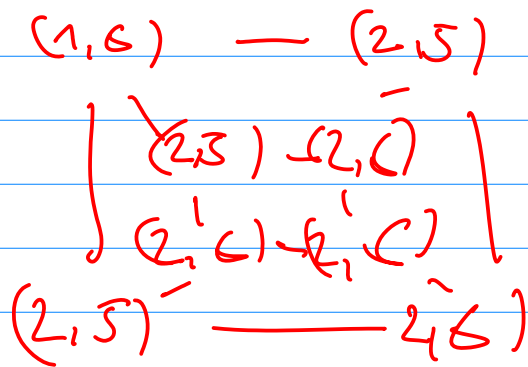
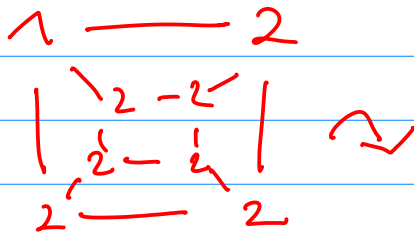
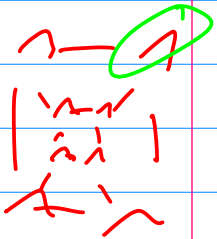
c) Scaling tree

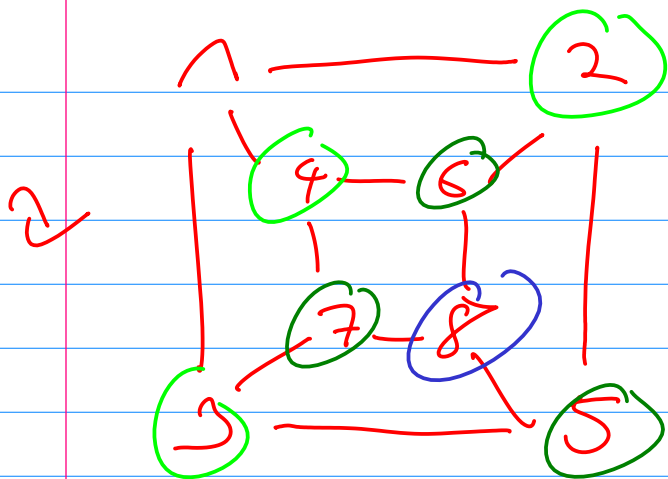


- choose an arbitrary scaling point

(note, in cube now the remaining 7 c are no longer identical)

- double marks and reduce the first by 1





routine CANON N atoms

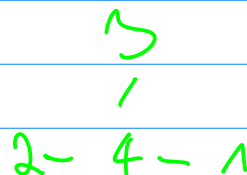
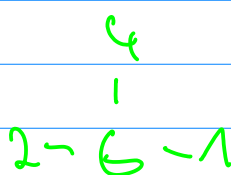
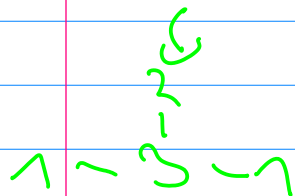
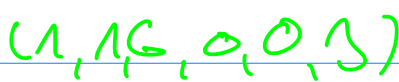
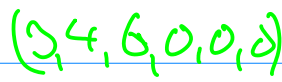
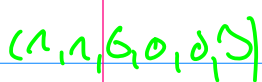
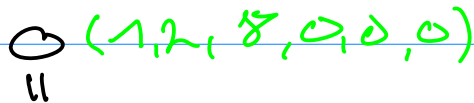
$$O(N^2 \log N)$$

GENES: generation of unique SMILES

1) initial node \rightarrow smallest canonical number

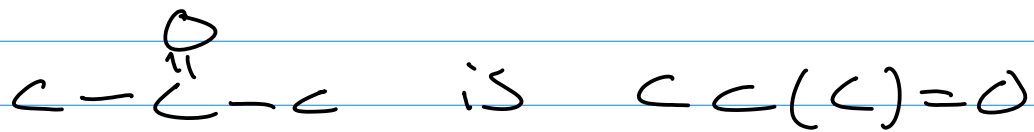
2) branching towards lower (at least) atom

Example: acetone



Therefore the canonical SMILES notation

for



and not CC(=O)C