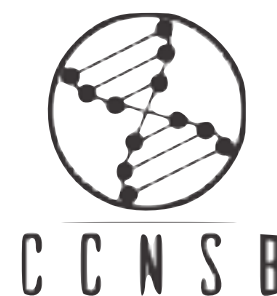


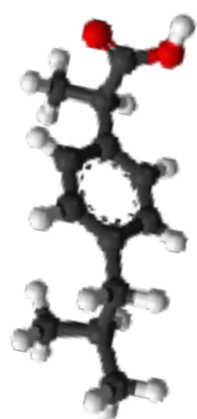
# Molecular Feature Vectors

(Representing molecules in a way that is good for ML)

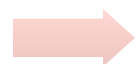


 @deva\_priyakumar

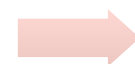
# Recap



$Z_i, R_i$



Software  
(eg. Gaussian16)



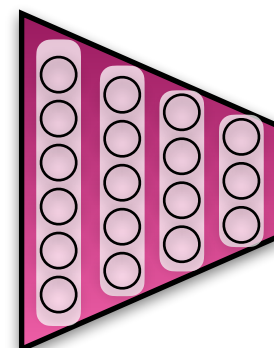
Properties  
(eg. Electronic Energy)

$$y = f(x)$$

$x$  and  $f$  are known;  
we compute  $y$



$x'$



$f$



Properties  
(eg. Solvation Free Energy)

$$y = f(x')$$

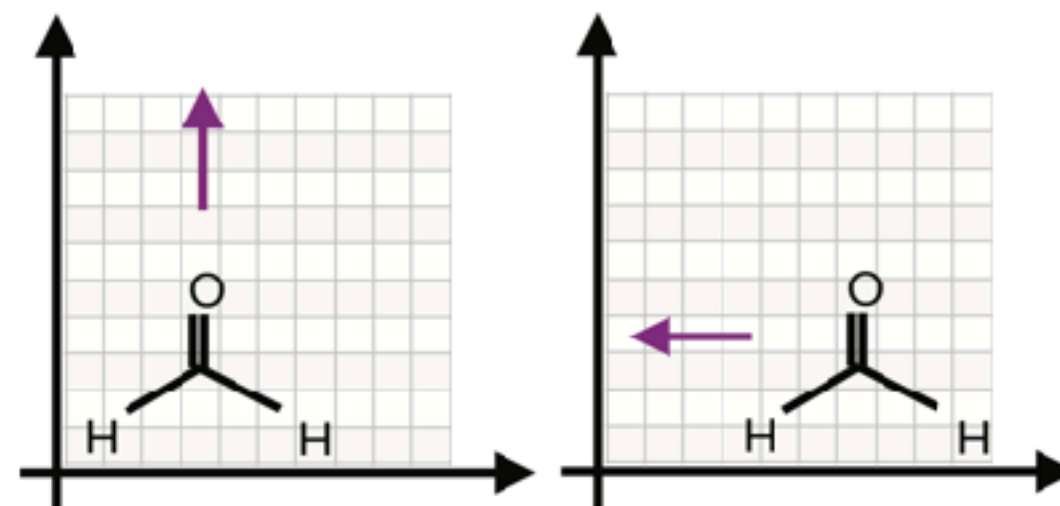
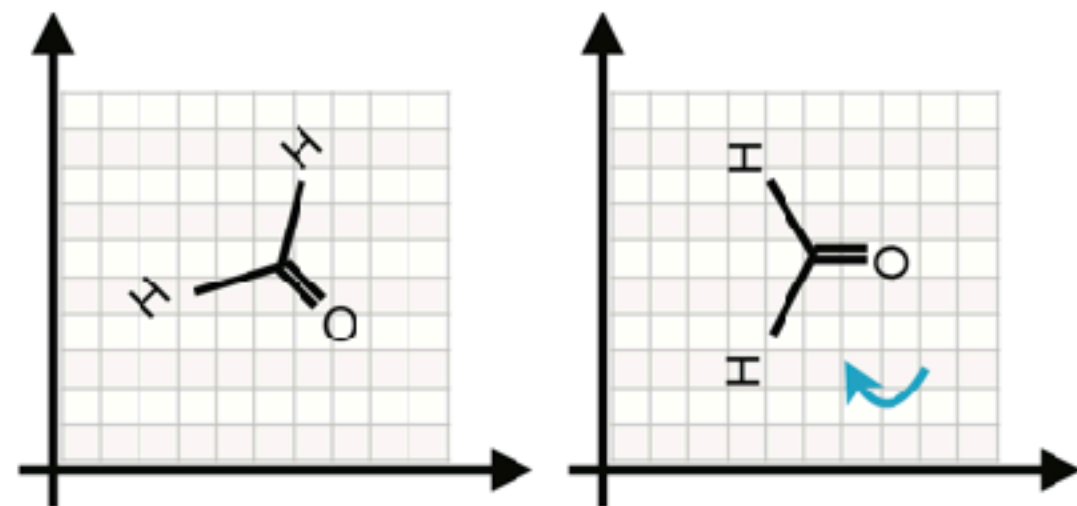
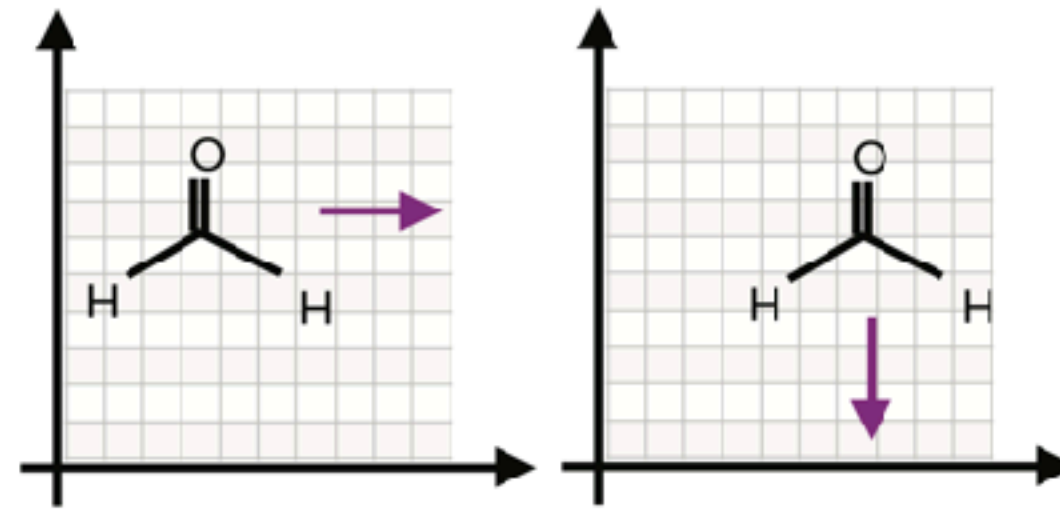
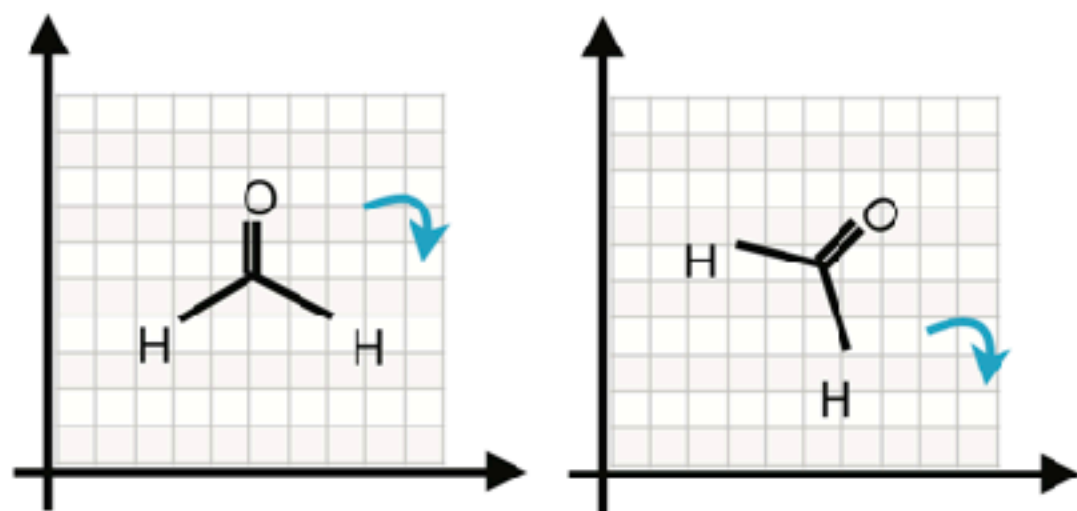
$x$  and  $y$  are known;  
we model  $f$

**Machine Learning**

**Blank**

**In a Computational Tool, How do we  
Represent Molecules?**

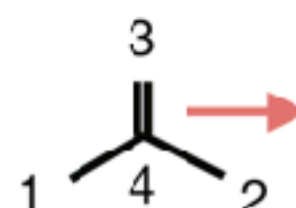
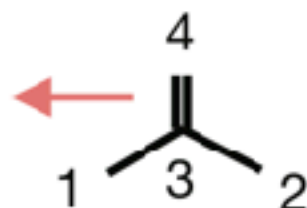
- Cartesian coordinates
- Internal coordinates
- Why unsuitable?
  - Not invariant to rotation
  - Not invariant to translation
  - Not invariant to change of atom ordering
  - Not a constant size feature vector



Rotations

Translations

$$\begin{bmatrix} \text{H} & 0 & 1 & 0 \\ 0 & \text{H} & 1 & 0 \\ 1 & 1 & \text{C} & 2 \\ 0 & 0 & 2 & \text{O} \end{bmatrix}$$

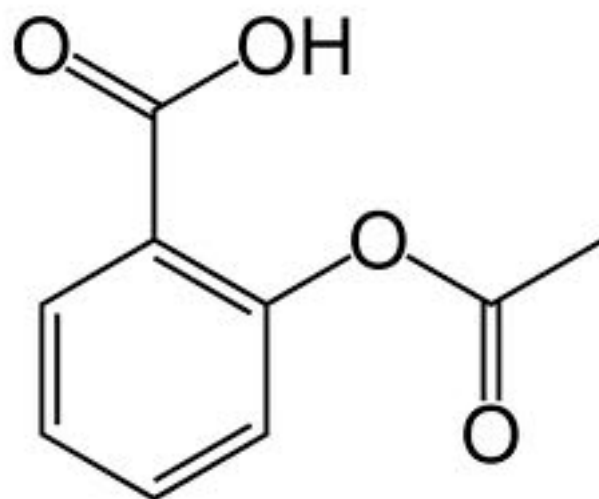


$$\begin{bmatrix} \text{H} & 0 & 0 & 1 \\ 0 & \text{H} & 0 & 1 \\ 0 & 0 & \text{O} & 2 \\ 1 & 1 & 2 & \text{C} \end{bmatrix}$$

Permutations

# Molecular Representation

- Name of the Molecule - Aspirin or 2-acetoxybenzoic acid
- Molecular Formula -  $C_9H_8O_4$
- Line diagram -



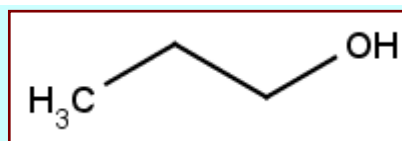




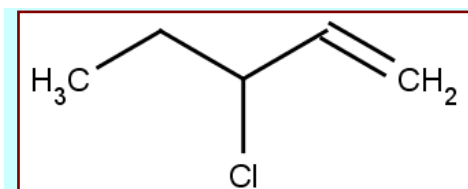
- Machine learning friendly featurizations
  - capture crucial information about the molecule
  - invariant to rotation
  - invariant to translation
  - invariant to the ordering of atoms
  - preferably a constant dimension for molecules of any size

# (1) SMILES

- Atoms given as atomic symbols
- No hydrogens (implicit)
- Double bonds represented as “=” and triple bonds by “#”
- Branches by parenthesis
- Rings represented by allocating digits to the connecting atoms

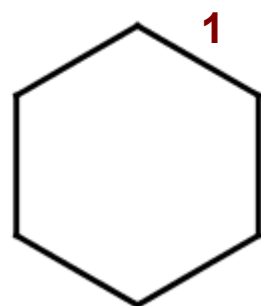


SMILES representation : **CCCCO**

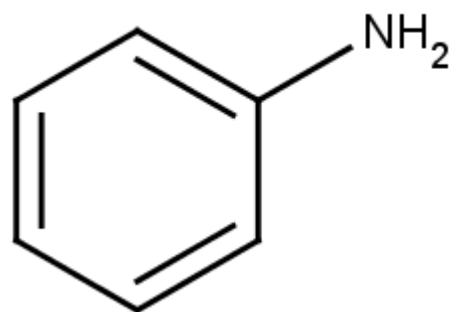


SMILES: **CCC(Cl)C=C**

- Rings represented by allocating digits to the connecting atoms
- Aromatic rings represented by lowercase

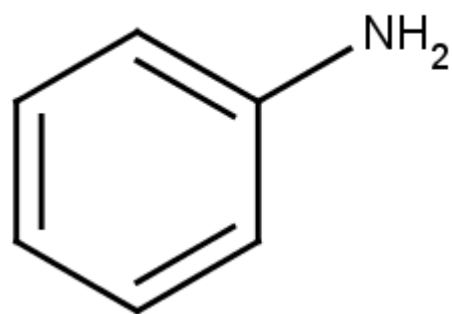


SMILES: **C1CCCCC1**



SMILES: **Nc1ccccc1**

- Unambiguous - a given SMILES string represents a unique structure.
- But, how about the reverse?



SMILES: **Nc1ccccc1**

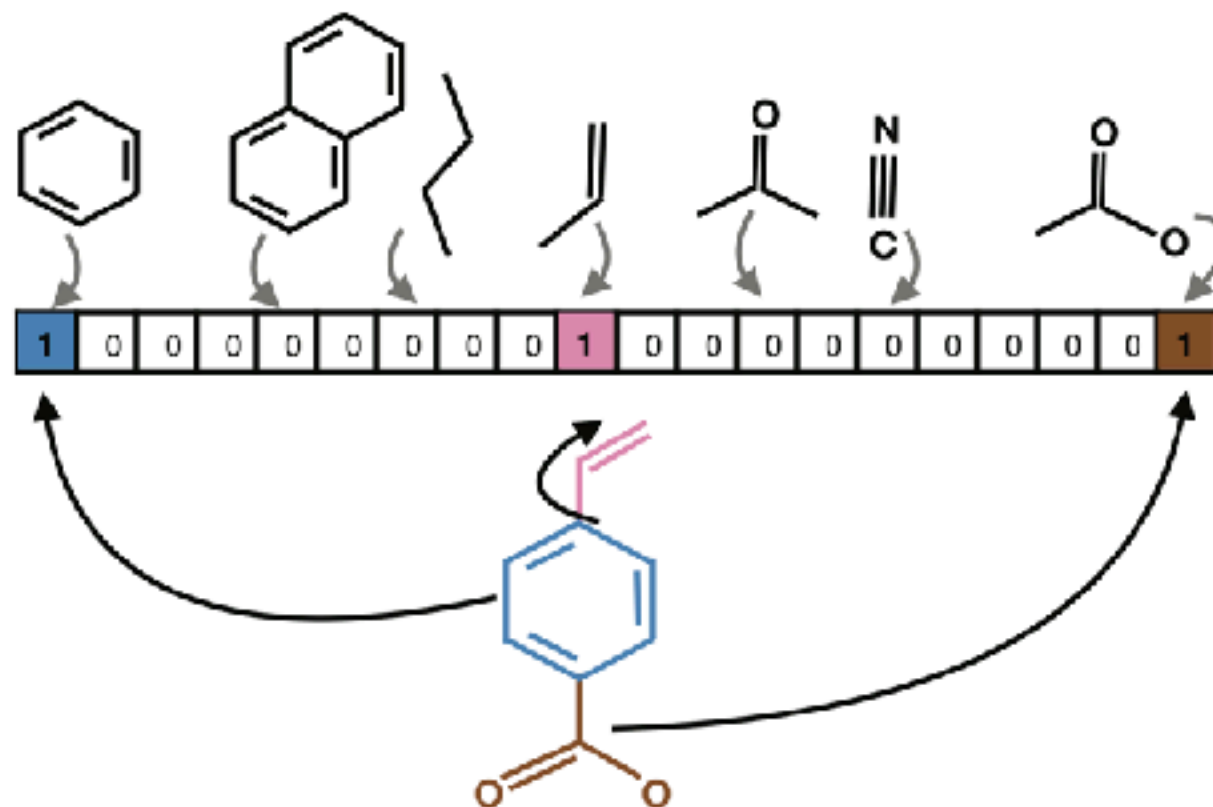
but also **c1ccccc1N**

or **c1cc(N)ccc1**

- Canonical representation of SMILES (**FIND OUT**)

## (2) Molecular Fingerprints

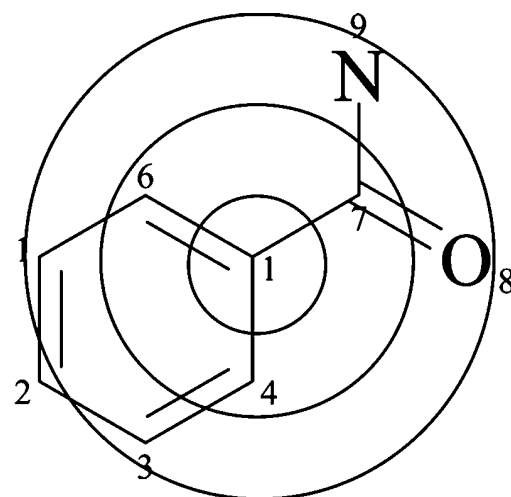
- Bit string - set the bits depending on the presence or absence of a given list of features
- Number of bits determined by the number of structural keys in the whole dataset.



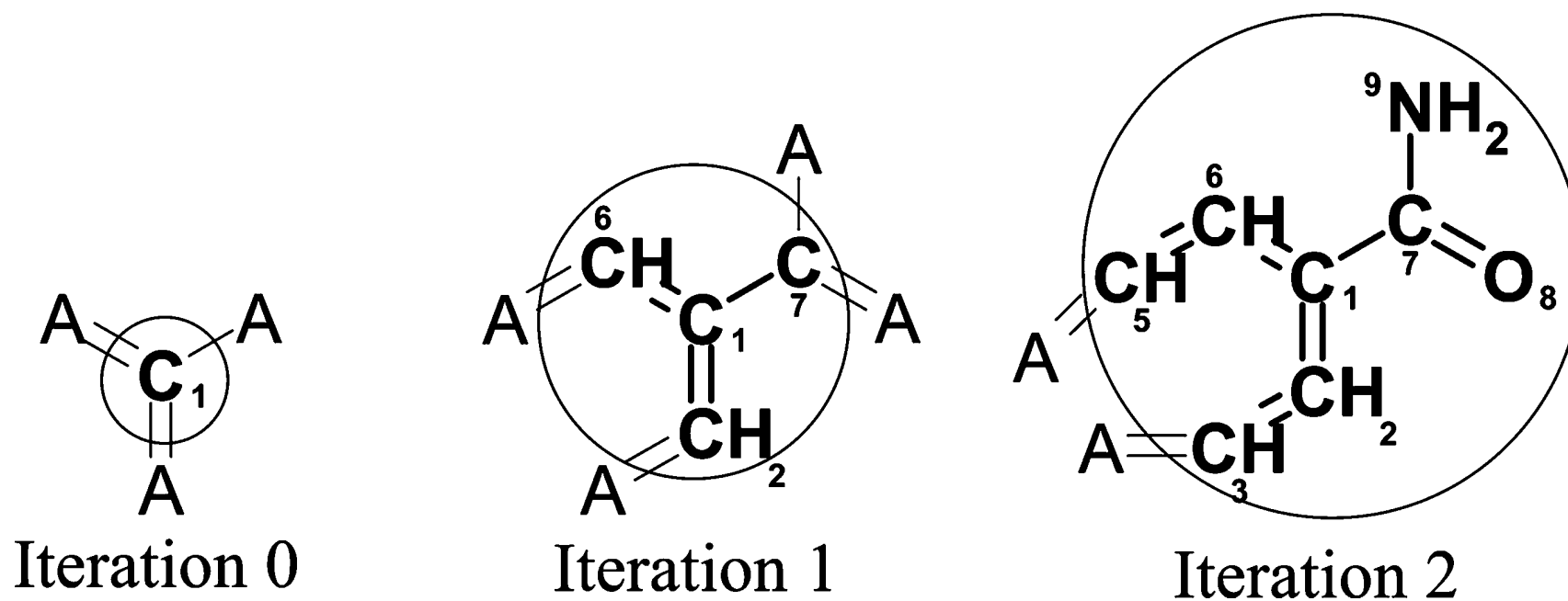
**FIGURE 3** Representation of a molecular fingerprint encoded the presence (1) or absence (0) of certain substructures in a compound. This molecule is represented by a vector of length 20 consisting of binary numbers

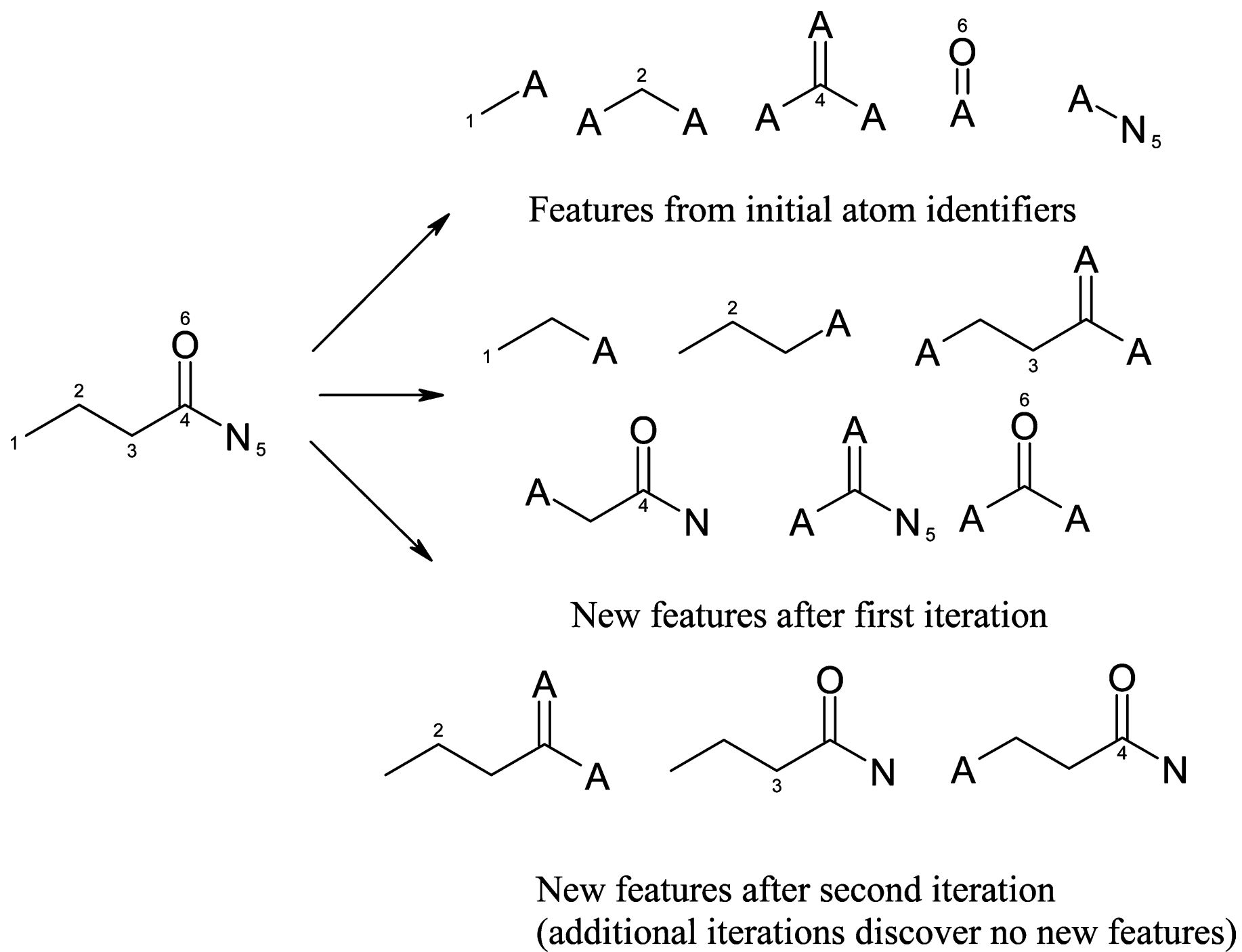
# Extended Connectivity Fingerprints (ECFP)

- Each molecule decomposed into submodules originating from each heavy atom
- Each of these assigned with a unique identifier
- This is extended through bonds to generate larger substructures (and corresponding identifiers)
- Hash all the substructures into a fixed length binary fingerprint representation.

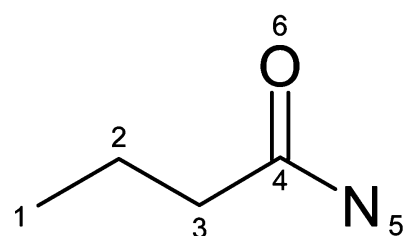


Considering atom 1 in benzoic acid amide









> <ECFP\_0>  
 734603939  
 1559650422  
 -1100000244  
 1572579716  
 -1074141656

> <ECFP\_2>  
 734603939  
 1559650422  
 -1100000244  
 1572579716  
 -1074141656  
 863188371  
 -1793471910  
 -1789102870  
 -1708545601  
 -932108170  
 2099970318

> <ECFP\_4>  
 734603939  
 1559650422  
 -1100000244  
 1572579716  
 -1074141656  
 863188371  
 -1793471910  
 -1789102870  
 -1708545601  
 -932108170  
 2099970318  
 -87618679  
 1112638790  
 -627599602

> <ECFP\_6>  
 734603939  
 1559650422  
 -1100000244  
 1572579716  
 -1074141656  
 863188371  
 -1793471910  
 -1789102870  
 -1708545601  
 -932108170  
 2099970318  
 -87618679  
 1112638790  
 -627599602

- Unique
- Rotational/translational invariant
- But, no information about the three-dimensional structure

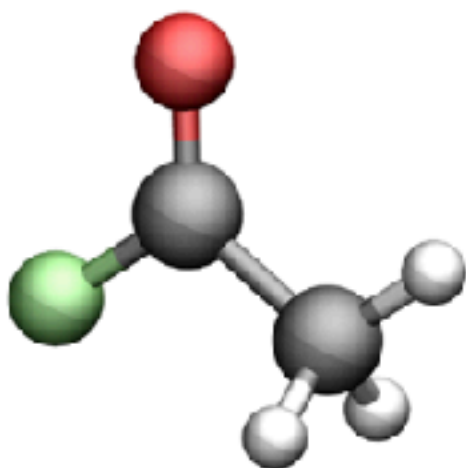
**Why do we need to represent the 3D  
structure?**

### (3) Coulomb Matrix

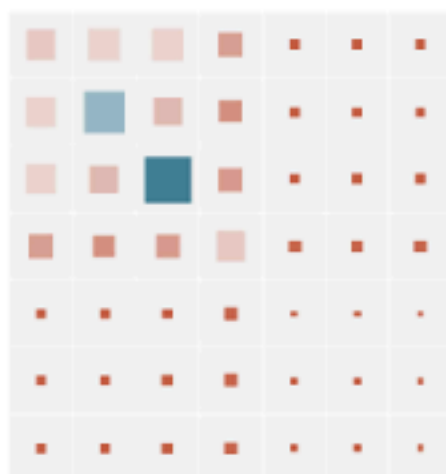
$$M_{IJ} = \begin{cases} 0.5Z_I^{2.4} & \text{for } I = J \\ \frac{Z_I Z_J}{|\mathbf{R}_I - \mathbf{R}_J|} & \text{for } I \neq J \end{cases}$$

- Off-diagonal elements => Coulomb repulsion between atoms I and J
- Diagonal elements => Polynomial fit of atomic self-energy

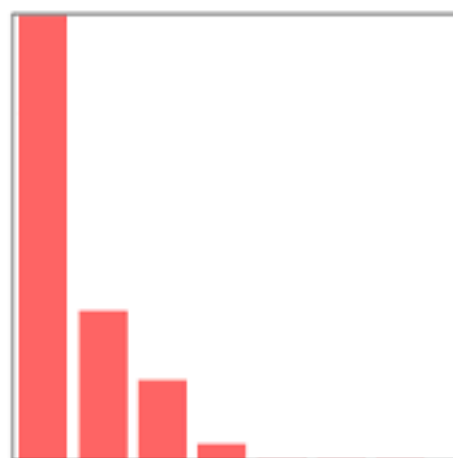
(A)



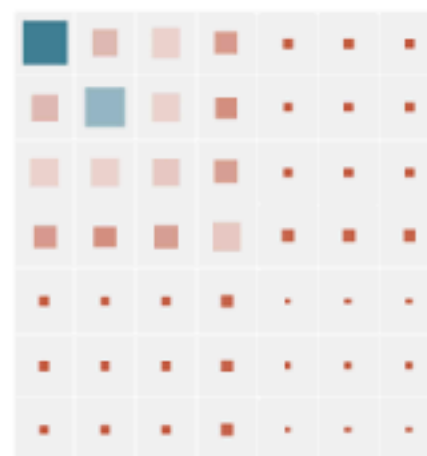
(B)



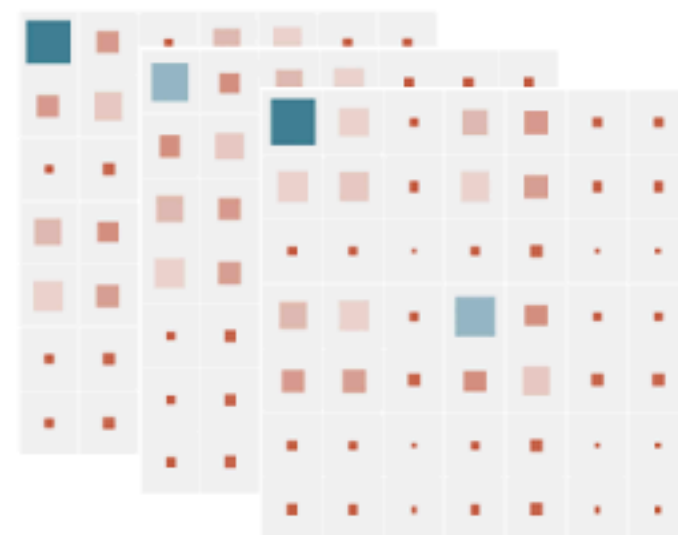
(C)



(D)

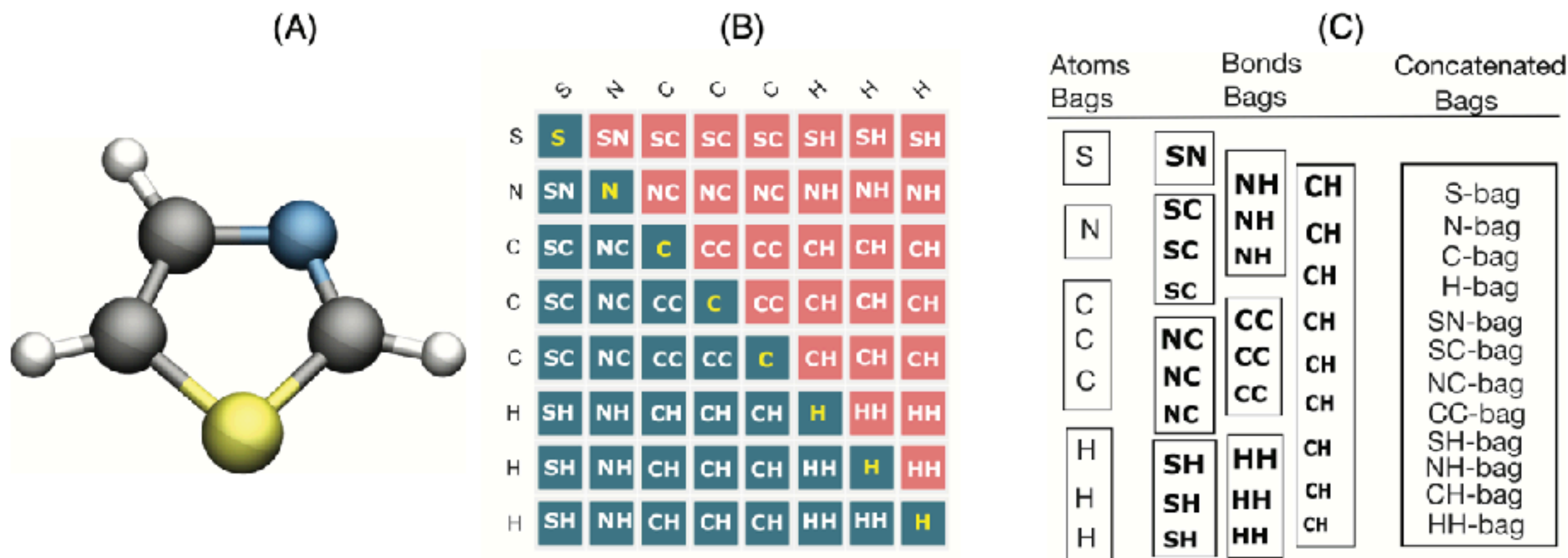


(E)



- Not invariant to swapping of atomic indices
- Tricks:
  - Use randomly permuted sets of matrices
  - Use eigen values of the Coulomb matrix...
- Not invariant to molecular size - pad with zeroes

## (4) Bag of Bonds



**FIGURE 5** Schematic view of the bag-of-bonds (BoB) representation: (A) ball and stick representation of the thiazole ( $C_3H_3NS$ ) molecule, (B) Coulomb matrix elements, (C) different Coulomb matrix entries (off-diagonal) are sorted into different bags. Here, concatenated bags are not padded with zeros for clarity

Each entry in the bag computed as

$$Z_i Z_j = \frac{Z_i Z_j}{|R_i - R_j|}$$

- Each entry in the bag is sorted
- Then sorted in a specific order
- Padded with zeros to account for molecular size invariance



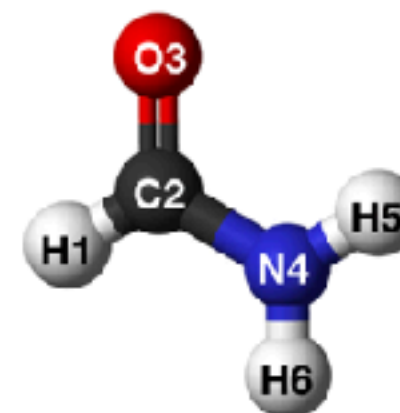
# (5) BAND

## (a) Atom identifier








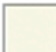










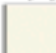




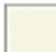









| H       | C       | O       | N       |
|---------|---------|---------|---------|
| 1 0 0 0 | 0 1 0 0 | 0 0 1 0 | 0 0 0 1 |

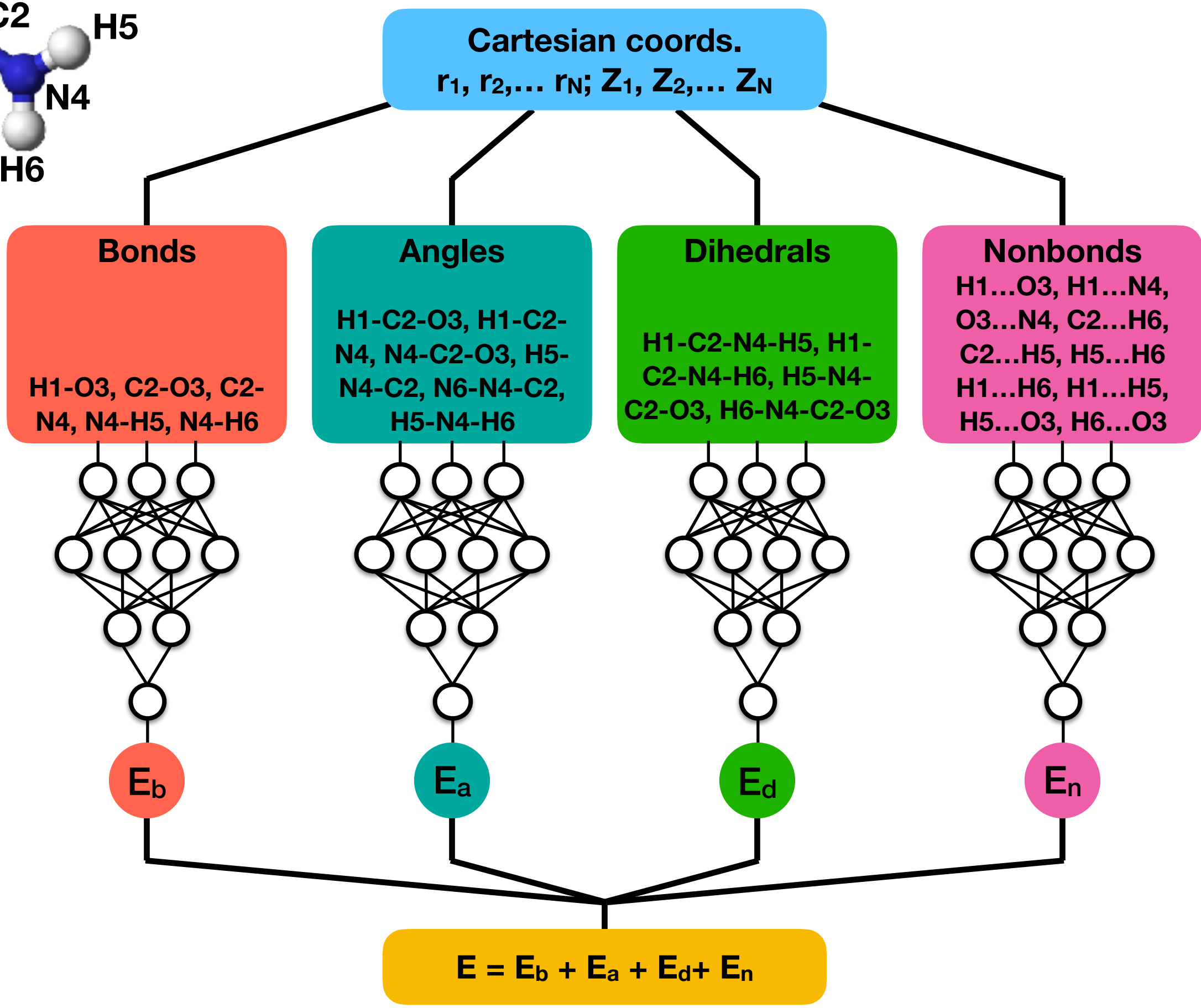
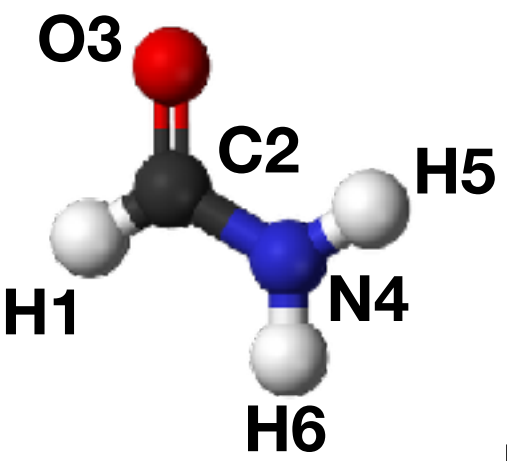
## (b) Atom identifier and atom typing

|    | Atom name | Atom type |
|----|-----------|-----------|
| H1 | 1 0 0 0   | 0 1 0 0   |
| N4 | 0 0 0 1   | 2 1 0 0   |



## (c) Feature vectors of bonds, angles, nonbonds and dihedrals

|          |  |  |   |  |  |   |   |  |  |  |               |
|----------|--|--|---|--|--|---|---|--|--|--|---------------|
| Bond     | atom <i>p</i><br>  | atom <i>q</i><br>  | <i>b<sub>pq</sub></i><br>   | 17 dimensions  |  |   |   |  |  |  |               |
| Angle    | atom <i>p</i><br>  | atom <i>q</i><br>  | atom <i>r</i><br>  | <i>a<sub>pqr</sub></i><br>  | <i>b<sub>pq</sub></i><br>   | <i>b<sub>qr</sub></i><br>  | 27 dimensions   |  |  |  |               |
| Nonbond  | atom <i>p</i><br>  | atom <i>q</i><br>  | <i>n<sub>pq</sub></i><br>   | 17 dimensions  |  |   |   |  |  |  |               |
| Dihedral | atom <i>p</i><br>  | atom <i>q</i><br>  | atom <i>r</i><br>  | atom <i>s</i><br>  | <i>d<sub>pqrs</sub></i><br> | <i>a<sub>pqr</sub></i><br> | <i>a<sub>qrs</sub></i><br> | <i>b<sub>pq</sub></i><br> | <i>b<sub>qr</sub></i><br> | <i>b<sub>rs</sub></i><br> | 38 dimensions |



**TO Learn:**

**SELFIES**

**ANI - Originally Symmetry Fn. by Behler & Parinello**

**SchNet**

# Slide added after the lecture

- <https://onlinelibrary.wiley.com/doi/full/10.1002/qua.26870>
- <https://jcheminf.biomedcentral.com/articles/10.1186/s13321-020-00460-5>