# Molecular Representation

- ***Keras*** provides a convenient programming interface for deep learning that allows us to have efficient/verified implementations of many standard procedures/algorithms relevant to deep neural networks
- For simple feed-forward, densely connected neural networks, one may use the **Sequential API**. For more complicated architectures, the **Functional API** may be needed.
- The basic procedure is *build* → *compile* → *fit* → *predict*
- In terms of optimizers, **NAdam** is usually a solid choice, but there are many options available
- In terms of loss functions, there are again many available, but it is easy to define your own should any of those be inadequate
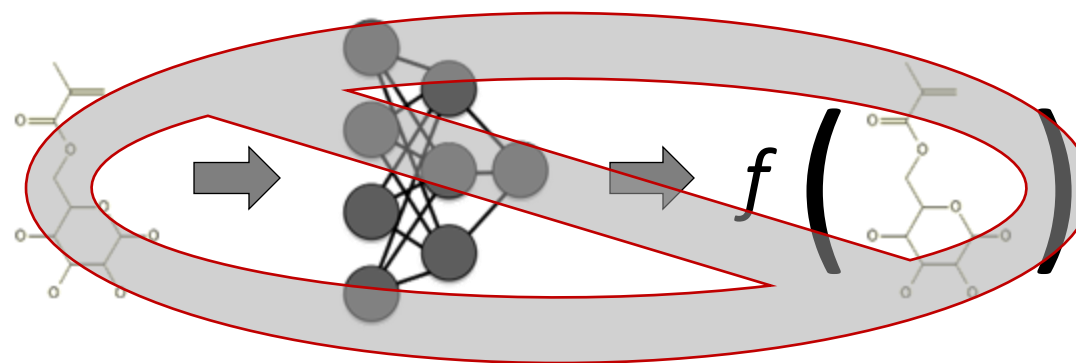
# Machine Learning Meets Molecules

A critical task for utilizing machine learning algorithms is _data representation_

_Consider the goal of developing Quantitative (Chemical) Structure Property Relationships...._



**Structure** ⇄ **Properties**
Prediction
Design

_we (as humans) know how to this process/goal, and we are interested in using ML to facilitate it_

_However, stick drawing or chemical name is a bit of problem for a neural network or other ML algorithm_

- Need technical methods to convert molecular structures into machine-readable formats (e.g., numerical vectors) that can be processed as inputs to ML algorithms
- The transformation should carry information that is useful to the prediction task in distinguishing amongst chemical moieties
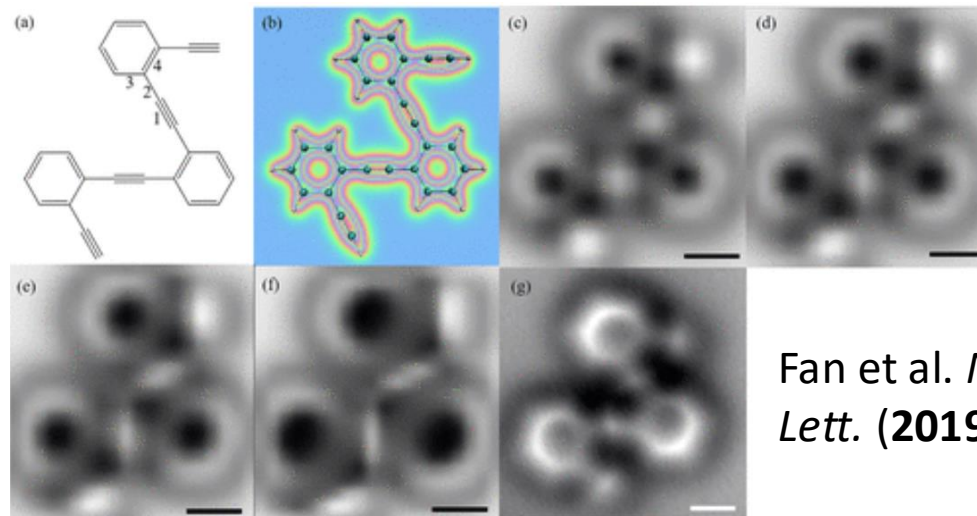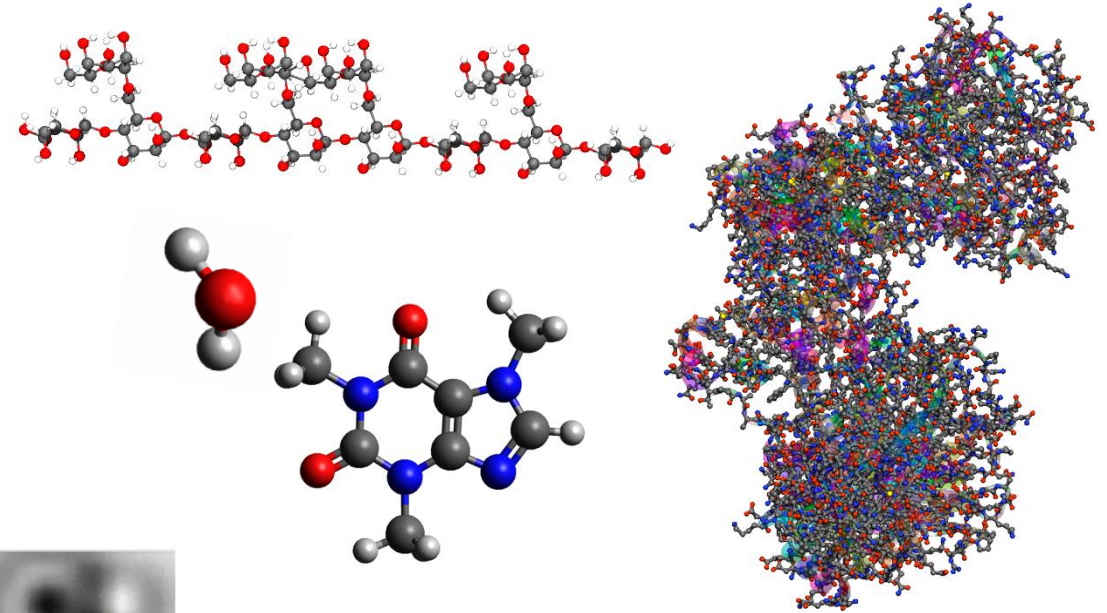
_This conversion process is referred to as **molecular featurization**_
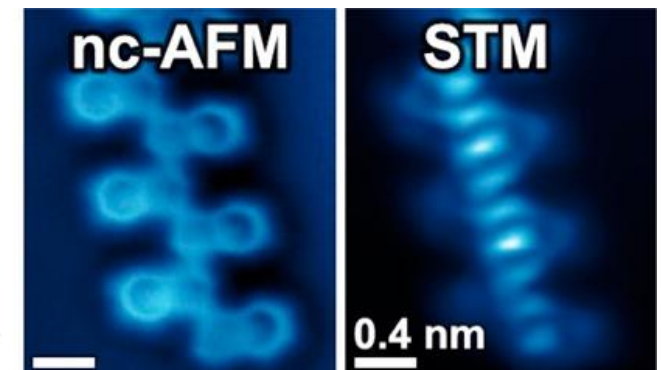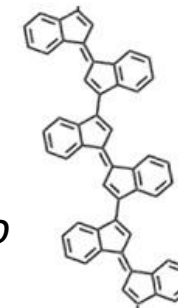
# What is a molecule?

Before we specifically address how to featurize a molecule, we should settle on what a molecule *is* since that may dictate the approach

## *Some possible characteristics of molecules*

- Groups of atoms that interact/are joined by physical forces?
- Notion of spatially localized electrons → bonds?
- Fundamental units of chemical reactions?
- Typified by specific chemical or physical properties?

Fan et al. *Nano Lett.* (**2019**)

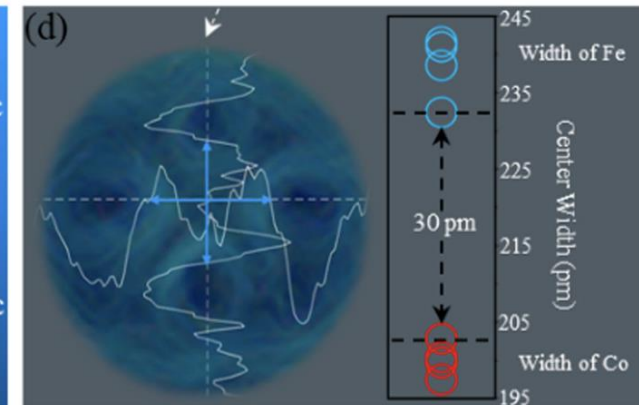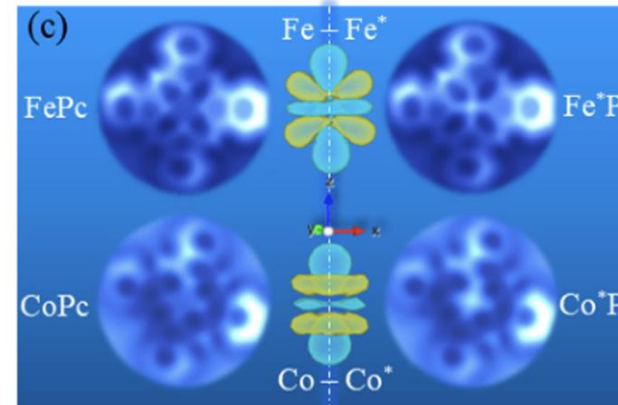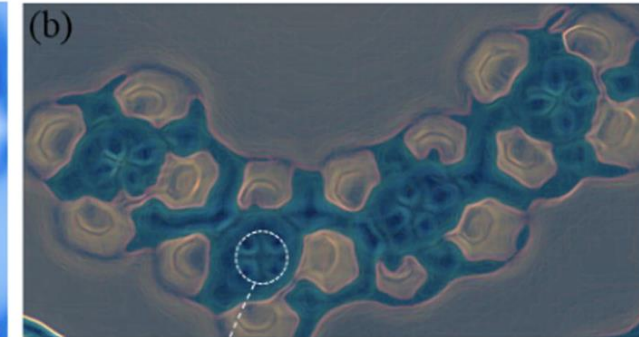Riss et al. *Nano Lett.* (**2014**)

# What is a molecule?

# Molecules as particles… or waves

```
 5 0  1
 6 C
 7 C  1  r2
 8 C  1  r3   2  a3
 9 C  2  r4   1  a4   3  d4
10 C  3  r5   1  a5   2  d5
11 C  5  r6   3  a6   1  d6
12 C  6  r7   5  a7   3  d7
13 C  7  r8   6  a8   5  d8
14 C  8  r9   7  a9   6  d9
15 H  1  r10  2  a10  3  d10
16 H  2  r11  1  a11  3  d11
17 C  8  r12  7  a12  6  d12
18 C  9  r13  8  a13  7  d13
19 C 12  r14  8  a14  7  d14
20 C 14  r15 12  a15  8  d15
21 H  3  r16  1  a16  2  d16
22 O  4  r17  2  a17  1  d17
23 O 12  r18  8  a18  7  d18
24 H 18  r19 12  a19  8  d19
25 O 13  r20  9  a20  8  d20
26 H 20  r21 13  a21  9  d21
27 C  5  r22  3  a22  1  d22
28 O 22  r23  5  a23  3  d23
29 O  7  r24  6  a24  5  d24
30 C 14  r25 12  a25  8  d25
31 C 25  r26 14  a26 12  d26
32 H 25  r27 14  a27 12  d27
33 H 25  r28 14  a28 12  d28
34 C 26  r29 25  a29 14  d29
35 C 15  r30 14  a30 12  d30
36 H 29  r31 26  a31 25  d31
37 H 29  r32 26  a32 25  d32
38 C 26  r33 25  a33 14  d33
39 C 33  r34 26  a34 25  d34
40 H 34  r35 33  a35 26  d35
41 H 34  r36 33  a36 26  d36
42 H 34  r37 33  a37 26  d37
43 O 26  r38 25  a38 14  d38
44 H 38  r39 26  a39 25  d39
45 O 33  r40 26  a40 25  d40
46 O 30  r41 15  a41 14  d41
47 H 30  r42 15  a42 14  d42
48 C 41  r43 30  a43 15  d43
49 C 43  r44 41  a44 30  d44
50 C 44  r45 43  a45 41  d45
51 C 45  r46 44  a46 43  d46
52 C 46  r47 44  a47 44  d47
53 H 47  r48 46  a48 45  d48
54 C 47  r49 46  a49 45  d49
55 C 44  r50 43  a50 41  d50
56 H 50  r51 44  a51 43  d51
57 H 50  r52 44  a52 43  d52
58 H 50  r53 44  a53 43  d53
59 O 45  r54 44  a54 43  d54
60 H 45  r55 44  a55 43  d55
61 H 54  r56 45  a56 44  d56
62 N 46  r57 45  a57 44  d57
63 H 46  r58 45  a58 44  d58
64 H 57  r59 46  a59 45  d59
65 H 57  r60 46  a60 45  d60
66 C 17  r61  4  a61  2  d61
67 H 61  r62 17  a62  4  d62
68 H 61  r63 17  a63  4  d63
69 H 61  r64 17  a64  4  d64
70 O 43  r65 41  a65 30  d65
71 H 43  r66 41  a66 30  d66
72 H 44  r67 43  a67 41  d67

73 Variables:
74 r2= 1.3887
75 r3= 1.3888
76 a3= 120.06
77 r4= 1.4109
78 a4= 120.62
79 d4=   0.30
80 r5= 1.4093
81 a5= 120.72
82 d5= 359.73
83 r6= 1.4027
84 a6= 119.93
85 d6= 359.70
86 r7= 1.4926
87 a7= 118.85
88 d7= 180.86
89 r8= 2.5252
90 a8=  92.09
91 d8= 359.58
92 r9= 1.4275
93 a9=  31.79
94 d9= 178.39
95 r10= 1.0827
96 a10= 119.98
97 d10= 179.86
98 r11= 1.0816
99 a11= 117.97
100 d11= 180.35
101 r12= 1.4298
102 a12= 150.94
103 d12= 179.57
104 r13= 1.4382
105 a13= 119.09
106 d13= 179.81
107 r14= 1.4038
108 a14= 120.67
109 d14= 358.38
110 r15= 1.4098
111 a15= 120.82
112 d15=   1.75
113 r16= 1.0825
114 a16= 118.14
115 d16= 179.73
116 r17= 1.3606
117 a17= 120.83
118 d17= 180.37
119 r18= 1.3583
120 a18= 123.77
121 d18= 179.16
122 r19= 0.9601
123 a19= 129.88
124 d19= 358.19
125 r20= 1.3555
126 a20= 123.41
127 d20= 179.19
128 r21= 0.9594
129 a21= 129.74
130 d21= 358.96
131 r22= 1.4792
132 a22= 119.70
133 d22= 179.59
134 r23= 1.2191
135 a23= 117.17
136 d23= 359.69
137 r24= 1.2178
138 a24= 117.65
139 d24= 180.74
140 r25= 1.5211
141 a25= 117.74
142 d25= 182.58
143 r26= 1.5473
144 a26= 114.42
145 d26= 161.99
146 r27= 1.1128
147 a27= 110.24
148 d27=  42.30
149 r28= 1.1126
150 a28= 106.57
151 d28= 285.79
152 r29= 1.5280
153 a29= 106.78
154 d29=  49.64
155 r30= 1.5397
156 a30= 120.60
157 d30= 178.95
158 r31= 1.1029
159 a31= 110.58
160 d31= 173.45
161 r32= 1.1141
162 a32= 108.83
163 d32=  55.68
164 r33= 1.5301
165 a33= 111.30
166 d33= 170.83
167 r34= 1.5038
168 a34= 120.95
169 d34= 178.31
170 r35= 1.1102
171 a35= 109.67
172 d35= 301.64
173 r36= 1.1109
174 a36= 110.58
175 d36= 182.16
176 r37= 1.1081
177 a37= 110.87
178 d37=  62.46
179 r38= 1.4093
180 a38= 110.26
181 d38= 291.74
182 r39= 0.9940
183 a39= 108.08
184 d39= 298.30
185 r40= 1.2234
186 a40= 120.56
187 d40=   0.30
188 r41= 1.4265
189 a41= 110.98
190 d41= 225.97
191 r42= 1.1146
192 a42= 104.41
193 d42= 104.16
194 r43= 1.4312
195 a43= 113.25
196 d43= 232.56
197 r44= 2.3580
198 a44=  92.98
199 d44= 182.11
200 r45= 1.5487
201 a45=  92.94
202 d45= 112.42
203 r46= 1.5404
204 a46= 108.47
205 d46= 331.08
206 r47= 1.5347
207 a47= 108.13
208 d47=  56.19
209 r48= 1.1055
210 a48= 111.45
211 d48= 179.37
212 r49= 1.1122
213 a49= 109.44
214 d49=  61.93
215 r50= 1.5306
216 a50= 142.31
217 d50= 242.48
218 r51= 1.1108
219 a51= 109.75
220 d51=  56.08
221 r52= 1.1117
222 a52= 110.75
223 d52= 296.69
224 r53= 1.1102
225 a53= 110.49
226 d53= 175.86
227 r54= 1.4067
228 a54= 109.63
229 d54=  91.63
230 r55= 1.1160
231 a55= 108.49
232 d55= 211.54
233 r56= 0.9935
234 a56= 107.06
235 d56= 157.24
236 r57= 1.4631
237 a57= 110.87
238 d57= 177.67
239 r58= 1.1150
240 a58= 109.22
241 d58= 297.59
242 r59= 1.0460
243 a59= 108.99
244 d59= 177.74
245 r60= 1.0470
246 a60= 108.94
247 d60=  61.04
248 r61= 1.3800
249 a61= 126.34
250 d61=   1.71
251 r62= 1.1092
252 a62= 109.08
253 d62= 179.30
254 r63= 1.1116
255 a63= 110.15
256 d63=  60.36
257 r64= 1.1106
258 a64= 110.34
259 d64= 298.54
260 r65= 1.4215
261 a65= 110.09
262 d65= 151.98
263 r66= 1.1134
264 a66= 111.17
265 d66=  36.18
266 r67= 1.1136
267 a67=  91.53
268 d67=   2.56
```
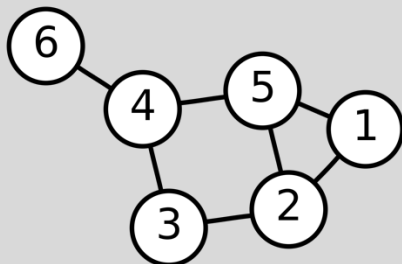
# Molecular Graphs

The idea of molecules being groups of connected atoms lends itself to representation as a *graph*



**Mathematical graph**

$$\mathcal{G} = (\mathcal{V}, \mathcal{E})$$

*Structure in discrete mathematics that usually demonstrates how some set of objects are related to one another*

$\mathcal{V}$ **Vertices/nodes** – indicates objects

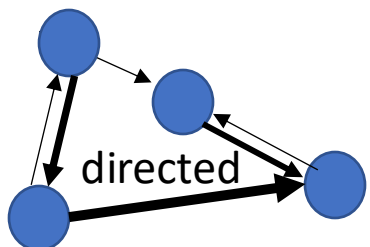$\mathcal{E}$ **Edges** – indicates pairwise relationship amongst objects

Graphs can be conveniently represented as matrices

$$\boldsymbol{G}, G_{ij} = \varepsilon_{ij}$$

**e.g.,**
$$\boldsymbol{G} = \begin{bmatrix} 0 & \varepsilon_{12} & 0 & 0 & \varepsilon_{15} & 0 \\ \varepsilon_{21} & 0 & \varepsilon_{23} & 0 & \varepsilon_{25} & 0 \\ 0 & \varepsilon_{32} & 0 & \varepsilon_{34} & 0 & 0 \\ 0 & 0 & \varepsilon_{43} & 0 & \varepsilon_{45} & \varepsilon_{46} \\ \varepsilon_{51} & \varepsilon_{52} & 0 & \varepsilon_{54} & 0 & 0 \\ 0 & 0 & 0 & \varepsilon_{64} & 0 & 0 \end{bmatrix}$$

graphs may be **directed** or **undirected**



directed

*vs.*

undirected

$$\varepsilon_{ij} = \varepsilon_{ji}$$

# Molecular Graphs

The idea of molecules being groups of connected atoms lends itself to representation as a *graph*

*In a molecular graph*



**Vertices** → Atoms or particles
**Edges** → bonds or interactions

$$\begin{pmatrix} 0 & 1 & 0 & 0 & 0 & 0 & 1 & 1 & 0 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \end{pmatrix}$$

***ignoring hydrogens**

Here, entries denote presence of a bond (or not) → often referred to as an **Adjacency matrix**, but the premise can easily be encode other bits of information

Nodes can report properties of each atom (element, charge, hybridization state)

Edges can indicate bond (order), distances, electronic properties, etc.

*This representation is typically deficient in the description of 3D structure/conformation, chirality, …*

# Text-based Representations

Common starting points for describing molecular graphs are text strings

**Such text strings should…**
- *be human-readable (not necessarily be intuitive)*
- *have well-defined rules to facilitate disambiguation*
- *Ideally possess canonicalization procedures*

*Most popular, pervasive:* **S**implified **M**olecular-**I**nput **L**ine **E**ntry **S**ystem (SMILES)

**Basic Rules**

1. *Atoms indicated by atomic symbols (<u>aromatic</u> rings → <u>lower case</u>)*
2. *<u>Inorganic</u> elements are enclosed by <u>brackets</u> (as are formal charges)*
3. *<u>Bonds</u> represented by <u>-,=,#, and :</u> (single, double, triple, and aromatic); single and aromatic bonds are conventionally omitted*
4. *<u>Branches</u> are specified by enclosures in <u>parentheses</u>*
5. *<u>Cyclic structures</u> are indicated by breaking one bond in each ring and <u>designating the point of opening/closure with a digit</u>*



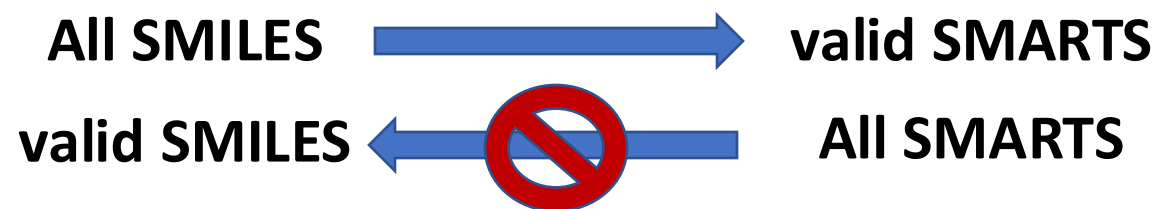*not natively canonical (additional algorithms do this)*

c1ccccc1C(=O)O

# Text-based Representations

Common starting points for describing molecular graphs are text strings

**An extension:** **SM**ILES **Ar**bitrary **T**arget **S**pecification (SMARTS)

**SMARTS** is not for representing molecular structures but *chemical patterns*

→ database queries
→ substructure searches
*(finding a subgraph of the molecular graph)*

**All SMILES** ➡ **valid SMARTS**

**valid SMILES** ⬅ 🚫 **All SMARTS**

**Decoding Exercise**

*draw out/describe the substructures from SMARTS*

- "cc"
- [c,n;H1]
- "Caa(O)aN"
- "Ca(aO)aaN"

## SMARTS Atomic Primitives

| Symbol | Symbol name | Atomic property requirements | Default |
|---|---|---|---|
| * | wildcard | any atom | (no default) |
| a | aromatic | aromatic | (no default) |
| A | aliphatic | aliphatic | (no default) |
| D<n> | degree | <n> explicit connections | exactly one |
| H<n> | total-H-count | <n> attached hydrogens | exactly one[1] |
| h<n> | implicit-H-count | <n> implicit hydrogens | at least one |
| R<n> | ring membership | in <n> SSSR rings | any ring atom |
| r<n> | ring size | in smallest SSSR ring of size <n> | any ring atom[2] |
| v<n> | valence | total bond order <n> | exactly one[2] |
| X<n> | connectivity | <n> total connections | exactly one[2] |
| x<n> | ring connectivity | <n> total ring connections | at least one[2] |
| -<n> | negative charge | -<n> charge | -1 charge (-- is -2, etc) |
| +<n> | positive charge | +<n> formal charge | +1 charge (++ is +2, etc) |
| #n | atomic number | atomic number <n> | (no default)[2] |
| @ | chirality | anticlockwise | anticlockwise, default class[2] |
| @@ | chirality | clockwise | clockwise, default class[2] |
| @<c><n> | chirality | chiral class <c> chirality <n> | (nodefault)[2] |
| @<c><n>? | chiral or unspec | chirality <c><n> or unspecified | (no default) |
| <n> | atomic mass | explicit atomic mass | unspecified mass |

## SMARTS Bond Primitives

| Symbol | Atomic property requirements |
|---|---|
| - | single bond (aliphatic) |
| / | directional bond "up"[1] |
| \ | directional bond "down"[1] |
| /? | directional bond "up or unspecified" |
| \? | directional bond "down or unspecified" |
| = | double bond |
| # | triple bond |
| : | aromatic bond |
| ~ | any bond (wildcard) |
| @ | any ring bond[1] |

## SMARTS Logical Operators

| Symbol | Expression | Meaning |
|---|---|---|
| exclamation | !e1 | not e1 |
| ampersand | e1&e2 | a1 and e2 (high precedence) |
| comma | e1,e2 | e1 or e2 |
| semicolon | e1;e2 | a1 and e2 (low precedence) |

# Text-based Representations

Common starting points for describing molecular graphs are text strings

**An extension: SM**ILES **Ar**bitrary **T**arget **S**pecification (SMARTS)
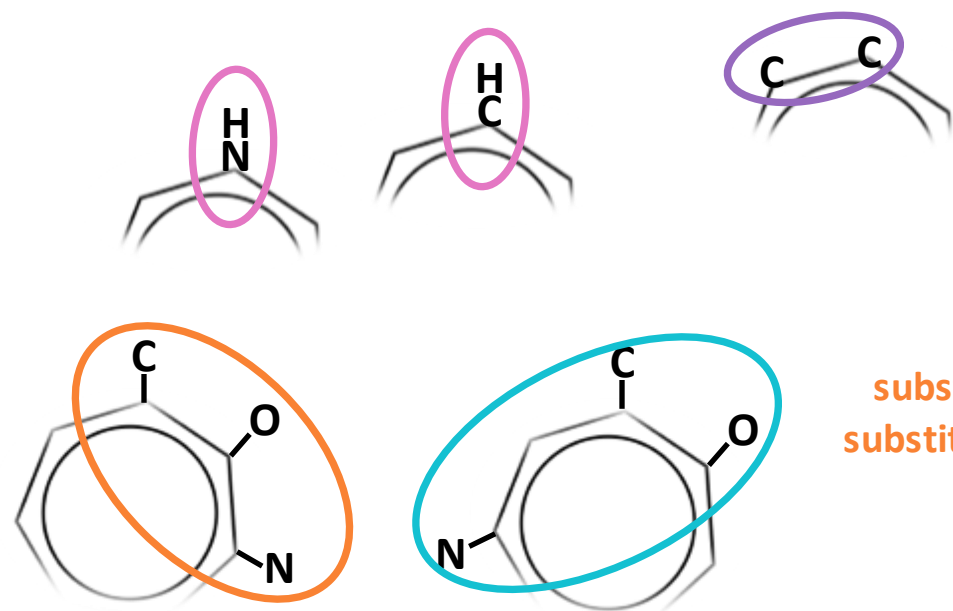
**SMARTS** is not for representing molecular structures but *chemical patterns* → database queries
→ substructure searches
*(finding a subgraph of the molecular graph)*

**All SMILES** → **valid SMARTS**

**valid SMILES** ← **All SMARTS**

**Decoding Exercise**

*draw out/describe the substructures from SMARTS*



any pair of bonded aromatic carbons
- "cc"

either aromatic carbon or nitrogen and exactly one hydrogen
- [c,n;H1]

substituent carbon of aromatic ring that is ortho to substituent oxygen and meta to substituent nitrogen
- "Caa(O)aN"

same as above but O and N likely para
- "Ca(aO)aaN"

# Text-based Representations

Common starting points for describing molecular graphs are text strings

*DeepSMILES* was developed to address some syntactic issues in using SMILES for "generative" models. The gist of the problem is that many perturbations to SMILES strings do not result in valid molecules.

| SMILES | DeepSMILES |
|---|---|
| C1CCCC1 | CCCCC5 |
| C1CCCCCCCC1 | CCCCCCCCC%10 |
| C(O)C | CO)C |
| C(OF)C | COF))C |
| C(F)(F)C | CF)F)C |
| C(=O)Cl | C=O)Cl |
| C(OC(=O)Cl)I | COC=O)Cl)))I |
| C1CC(OC)CC1 | CCCOC))CC5 |
| C1=C/CCCCC/1 | C=C/CCCCC/8 |
| C\1=C/CCCCC1 | C=C/CCCCC/8 |
| B(c1ccccc1)(O)O | Bcccccc6))))))O)O |
| Cn1cccc-2nccc12 | Cncccccnccc9-5 |
| C1N[C@@]12CO2 | CN[C@@]3CO3 |
| [C@@]12(NC1)CO2 | [C@@]NC3))CO3 |
| CC1CCCO[C@]21CCCCO2 | CCCCCO[C@@]6CCCCO6 |
| CC1CCCO[C@@]12CCCCO2 | CCCCCO[C@@]6CCCCO6 |
| NC[C@]12CCCC1C3CC2CC3 | NC[C@]CCCC5CCC8CC5 |
| NC[C@]12CCCC2C3CC1CC3 | NC[C@@]CCCC5CCC8CC5 |

**Other string-based representations**
*Wiswesser line notation*
*SYBYL Line Notation*
*IUPAC **In**ternational **Ch**emical **I**dentifier*

# Text-based Representations

Common starting points for describing molecular graphs are text strings

**Self**-referencing **E**mbedded **S**trings (SELFIES)

- New kid on the block with growing utility
- developed as a "100% robust" alternative to SMILES:
  - *every SELFIES string is a valid molecule*
  - *every molecule has a SELFIES*

# Text-based Representations

Common starting points for describing molecular graphs are text strings

**Self**-referencing **E**mbedded **S**trings (SELFIES)

**Formal Grammar Rules**



*conversion to molecular graph*

[F][=C][=C][#N]

1. Start in $X_0$ → F $X_1$
2. F $X_1$ → F C $X_3$
3. F C $X_3$ → F C = C $X_2$
4. F C = C $X_2$ → F C = C = N