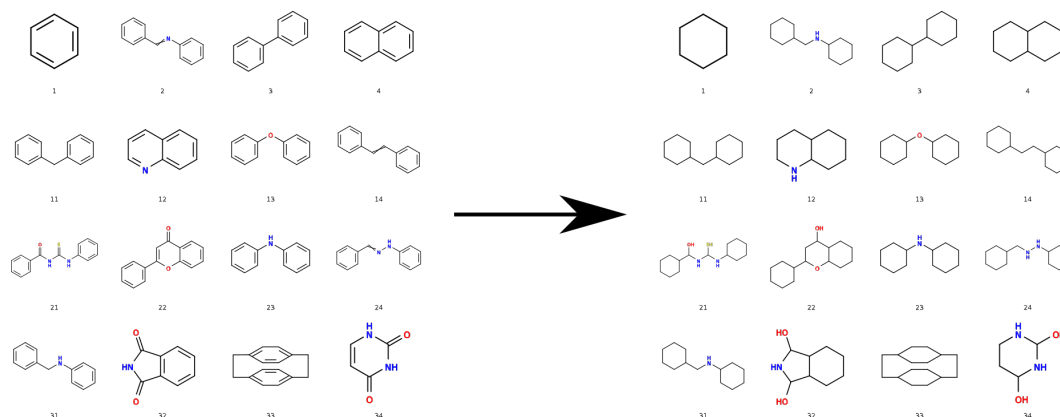


1 Background

The Bemis-Murcko scaffold¹ provided by DataWarrior² retains information about bond order and chirality. Sometimes, however, it suffices to retain only atom connectivity, like an assumption «there are only single bonds». Note DataWarrior equally offers the export of Bemis-Murcko skeleton, however this simplifies e.g. the scaffold about an imidazole into one of cyclopentane.



2 Typical use

The script runs from Python's CLI with a file listing SMILES to process as parameter. File `test_input.smi` (from sub-folder `test_data`) is an example:

```
python saturate_murcko_scaffolds.py [test_input.smi]
```

This generates `test_input_sat.smi` as permanent record; the addition of `_sat` is only a reminder of the performed saturation. The input file is preserved.

The file extension `.smi` of the input file is a suggestion, because it is frequently seen (e.g., around OpenBabel³). Internally, the script considers any character prior to the first period as part of the name of the input file. The name of the report file is a mere concatenation of this and the string `_sat.smi`.

While written for current branch of Python 3.6+ (e.g., Python 3.9.1), the script equally works well with the legacy branch of Python 2 like Python 2.7.16.

¹Bemis GW, Murcko MA *J. Med. Chem.* 1996, **39**, 2887-2893, doi 10.1021/jm9602928.

²Sander T, Freyss J, von Korff M, Rufener C, *J. Chem. Inf. Model.* 2015, **55**, 460-473, doi 10.1021/ci500588j. The program, (c) 2002–2021 by Idorsia Pharmaceuticals Ltd., is freely available under <http://www.openmolecules.org>. For the source code (GPLv3), see <https://github.com/thsa/datawarrior>.

³www.openbabel.org. The script initially was developed for and tested with OpenBabel (release 2.4.1; Nov 12, 2018) and Python 2.7.17 provided by Linux Xubuntu 18.04.2 LTS. It equally works with Python 3.9.1+ (released January 20, 2021) and OpenBabel (release 3.1.1 by January 6, 2021) as provided in Debian 10.

3 Example

For a collection of organic materials, the Bemis-Murcko scaffolds were extracted with DataWarrior (then release 5.0.0 for Linux, January 2019) as listing `test_input.smi` including higher bond orders (see folder `test_data`). The effect of the «artificial saturation» is easy to recognize while comparing the scaffold lists (fig. 1) in a difference view of the two `.smi` files.

013	<chem>c(cc1)ccc10c1cccc1</chem>	013	<chem>C(CC1)CCC10C1CCCC1</chem>
014	<chem>C(c1cccc1)=C/c1cccc1</chem>	014	<chem>C(C1CCCC1)CC1CCCC1</chem>
015	<chem>c1cc2cc3cccc3cc2cc1</chem>	015	<chem>C1CC2CC3CCCC3CC2CC1</chem>
016	<chem>O=C(c1cccc1)c1cccc1</chem>	016	<chem>OC(C1CCCC1)C1CCCC1</chem>
017	<chem>c1c[nH]c2c1cccc2</chem>	017	<chem>C1C[NH]C2C1CCCC2</chem>
018	<chem>c(cc1)ccc1/N=N/c1cccc1</chem>	018	<chem>C(CC1)CCC1NNC1CCCC1</chem>
019	<chem>C(c1cccc1)=N/N=C/c1cccc1</chem>	019	<chem>C(C1CCCC1)NNC1CCCC1</chem>
020	<chem>C(Cc1cccc1)c1cccc1</chem>	020	<chem>C(CC1CCCC1)C1CCCC1</chem>
021	<chem>O=C(c1cccc1)NC(Nc1cccc1)=S</chem>	021	<chem>OC(C1CCCC1)NC(NC1CCCC1)S</chem>
022	<chem>O=C1c(cccc2)c2OC(c2cccc2)=C1</chem>	022	<chem>OC1C(CCCC2)C2OC(C2CCCC2)C1</chem>
023	<chem>c(cc1)ccc1Nc1cccc1</chem>	023	<chem>C(CC1)CCC1NC1CCCC1</chem>
024	<chem>C(c1cccc1)=N/Nc1cccc1</chem>	024	<chem>C(C1CCCC1)NNC1CCCC1</chem>
025	<chem>O=C(C=CN1[C@@H]2OCCC2)NC1=O</chem>	025	<chem>OC(CCN1[C@@H]2OCCC2)NC1O</chem>
026	<chem>c1ccc2c(-c3cccc4cccc34)cccc2c1</chem>	026	<chem>C1CCC2C(-C3CCCC4CCCC34)CCCC2C1</chem>
027	<chem>c1ccc(C(c2cccc2)c2cccc2)cc1</chem>	027	<chem>C1CCC(C(C2CCCC2)C2CCCC2)CC1</chem>
028	<chem>c(cc1)cc2c1[nH]c1c2cccc1</chem>	028	<chem>C(CC1)CC2C1[NH]C1C2CCCC1</chem>
029	<chem>c(cc1)ccc1P(c1cccc1)c1cccc1</chem>	029	<chem>C(CC1)CCC1P(C1CCCC1)C1CCCC1</chem>
030	<chem>c1c(-c2cccc2)oc2c1cccc2</chem>	030	<chem>C1C(-C2CCCC2)OC2C1CCCC2</chem>
031	<chem>C(c1cccc1)Nc1cccc1</chem>	031	<chem>C(C1CCCC1)NC1CCCC1</chem>
032	<chem>O=C(c1c2cccc1)NC2=O</chem>	032	<chem>OC(C1C2CCCC1)NC2O</chem>
033	<chem>C(Cc1ccc(CC2)cc1)c1ccc2cc1</chem>	033	<chem>C(CC1CCC(CC2)CC1)C1CCC2CC1</chem>

Figure 1: Difference view of the SMILES strings of a Murcko scaffold *prior* (left hand column) and *after* an «artificial saturation» (right hand column). The processing affects explicit bond order indicators, e.g. double bond (equality sign, e.g., line #14), triple bond bond (octohorpe, not shown); or about implicit aromatization (lower case → upper case) for atoms of carbon, nitrogen, oxygen (depicted); or phosphorus, sulfur (not depicted). Stereochemical indicators about double bonds will be removed (e.g., slashes in lines #18 and #19). Descriptors of stereogenic centers (@-signs, e.g., line #25) are copied verbatim.

Subsequently, OpenBabel³ was used to illustrate the work performed. While eventually automated (cf. script `test_series.py`, deposit in folder `test_data`), instructions issued to OpenBabel on the command line followed the pattern of

```
2 obabel -ismi test_input.smi -O test_input_color.svg -xc10 -xr12 -xl --addinindex
```

to generate a `.svg` file (vector representation), or

```
3 obabel -ismi test_input_sat.smi -O test_input_sat_color.png -xc10 -xr12 -xl  
↪ --addinindex -xp 3000
```

to generate a bitmap `.png` with structure formulae depicted in a grid of 10 columns by 12 rows.

It is remarkable how well OpenBabel's displays the molecular structures with advanced motifs. In addition to those shown in the first illustration of this guide, see sub-folder `test_data` for a more extensive survey (e.g., the scaffold of cyclophane [entry #33], sparteine [#38], or adamantane [#50]).

4 Known peculiarities

The script neither removes, nor newly assigns SMILES descriptors about the absolute configuration of stereogenic centers (@). Thus, the «reduction» of double bonds e.g., ketones to secondary alcohols may yield new stereogenic centers with a description incomplete in this regard.

To resolve implicitly described aromatic systems, the script capitalizes the characters c, n, o, p, and s about the elements more frequently involved in ring formation. To avoid ambiguity processing data about tin – typically described by SMILES strings as [Sn], which the algorithm would convert into [SN] – *any* SMILES string including either the pattern of [Sn] or [sn] is excluded from the reduction. This rule is applied applied to both tin analogues of benzene, e.g. c1[sn]cccc1, and SMILES strings describing tin in a side chain (e.g., about a Stille reagent). Script `saturate_murcko_scaffolds.py` annotates these SMILES strings accordingly in the output file.

The script will not actively alter a charge assigned to an atom. If present (e.g., quaternary ammonium, carboxylate), this information will be carried over to the newly written SMILES string. Given the reduction of bond orders, depending on the substrate submitted, this approach may be sensible (e.g., about N in cetyltrimethylammonium bromide), or not (e.g., about N in pyridine *N*-oxide). Other libraries than the current script (e.g., RDKit⁴) might offer help to sanitize the processed SMILES strings.

If the input SMILES string describes more than exactly one molecule by the concatenating "." (period character), this special sign equally is the newly written SMILES string. This permits working with SMILES about e.g., co-crystals, like about 1,4-benzoquinone and hydroquinone, C1=CC(=O)C=CC1=O.c1cc(ccc1O)O.

5 License

Norwid Behrnd, 2019–21, GPLv3.

⁴For an overview about the freely available RDKit library, see www.rdkit.org. An introduction into the topic of «molecular sanitization» is provided in the section of this very title in the on-line [RDKit Book](#).