

A stub a day keeps the docstrings at bay

Paolo Tosco 12th RDKit UGM, Mainz September 21, 2023

YYYYXYYYYY



Current RDKit Python docstring status

- No hints at all in Visual Studio Code
 - VSCode uses Pyright for hints
 - Pyright needs stubs, and RDKit has none

```
from rdkit import Chem
mol = Chem.MolFromSmiles("CCC")
            (function) GetAtomWithIdx: Any
atom = mol.GetAtomWithIdx(
```

```
mol = Chem.MolFromSmiles("CCC")
Last executed at 2023-08-04 12:30:55 in 8ms
mol.GetAtomWithIdx(
                     Docstring:
                     GetAtomWithIdx( (Mol)arg1, (int)arg2) -> Atom :
                         Returns a particular Atom.
                           ARGUMENTS:
```

- Some hints in Jupyter Lab, however
 - arg1 should rather read self
 - arg2 should rather read idx



- idx: which Atom to return

NOTE: atom indices start at 0

gen_rdkit_stubs.py

- Based on <u>pybind11_stubgen</u>
 - Implements a docstring preprocessing hook to edit boost::python docstrings before feeding them to pybind11_stubgen
 - Traverses the rdkit module to generate all relevant stubs
 - Many non-static methods are incorrectly labelled as staticmethod due to lack of self parameter in the docstring signature
 - Any return types

```
from rdkit import Chem
mol = Chem.MolFromSmiles("CCC")
 (variable) atom: Any
atom = mol.GetAtomWithIdx(0)
             (method) def GetAtomWithIdx(
                 arg1: Mol.
                 arg2: int
             ) -> Atom
             GetAtomWithIdx( arg1: Mol, arg2: int) -> Atom
               Returns a particular Atom.
                ARGUMENTS: - idx: which Atom to return
                NOTE: atom indices start at 0
               C++ signature:
                 RDKit::Atom* GetAtomWithIdx(RDKit::ROMol {Ivalue},unsigned int)
       (function) GetFormalCharge: Any
atom.GetFormalCharge()
```

patch_rdkit_docstrings

- Based on clang AST parsing
 - Runs in parallel through multiprocessing
 - Generates AST file for each C++
 RDKit Python wrapper
 - Finds docstrings that need missing self parameter
 - Finds docstrings with arg1, arg2, ...
 parameter names and replaces them with parameter names extracted from C++ function signatures
 - Patches C++ sources

```
from rdkit import Chem
mol = Chem.MolFromSmiles("CCC")
 (variable) atom: Any
atom = mol.GetAtomWithIdx(0)
             (method) def GetAtomWithIdx(
                  arg1: Mol,
                 arg2: int
              ) -> Atom
             GetAtomWithIdx( arg1: Mol, arg2: int) -> Atom
               Returns a particular Atom.
                ARGUMENTS: - idx: which Atom to return
                NOTE: atom indices start at 0
               C++ signature:
                  RDKit::Atom* GetAtomWithIdx(RDKit::ROMol {Ivalue},unsigned int)
       (function) GetFormalCharge: Any
atom.GetFormalCharge()
```

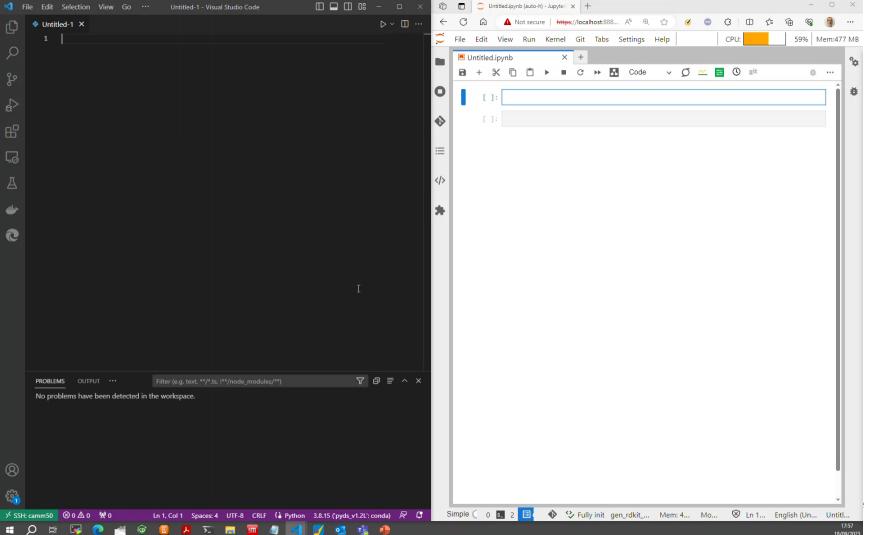
New RDKit D.O.C. strings

The extent of automated patching is quite large:

```
- $ git diff | grep -c '^+'
1741
$ git diff | grep -c '^-'
1180
```

 This is not something that can be done manually, unless one has a lot of time and patience on his/her hands

```
from rdkit import Chem
mol = Chem.MolFromSmiles("CCC")
 (variable) atom: Atom
atom = mol.GetAtomWithIdx(0)
              (method) def GetAtomWithIdx(idx: int) -> Atom
             GetAtomWithIdx( self: Mol, idx: int) -> Atom
                Returns a particular Atom.
                ARGUMENTS: - idx: which Atom to return
                NOTE: atom indices start at 0
                C++ signature:
                  RDKit::Atom* GetAtomWithIdx(RDKit::ROMol {Ivalue},unsigned int)
       (method) def GetFormalCharge() -> int
       GetFormalCharge( self: Atom) -> int
         C++ signature:
           int GetFormalCharge(RDKit::Atom {Ivalue})
atom.GetFormalCharge()
```



Thank you

YYYYXYYYYY

