

AMI2: High-throughput extraction of semantic chemistry from the scientific literature – the ChemistryVisitor

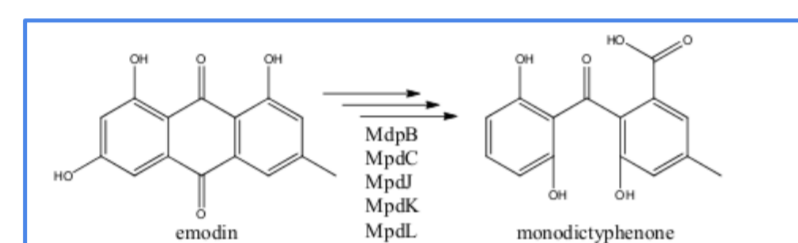
Andy Howlett, Mark Williamson, Peter Murray-Rust, Robert Glen
Unilever Centre, Cambridge

The AMI2 project aims to extract facts from the scientific literature, in particular from diagrams, which have been used less than text in the past. As part of this, the ChemistryVisitor is responsible for processing diagrams containing molecules and reactions / reaction schemes.

Metabolites 2012, 2 119

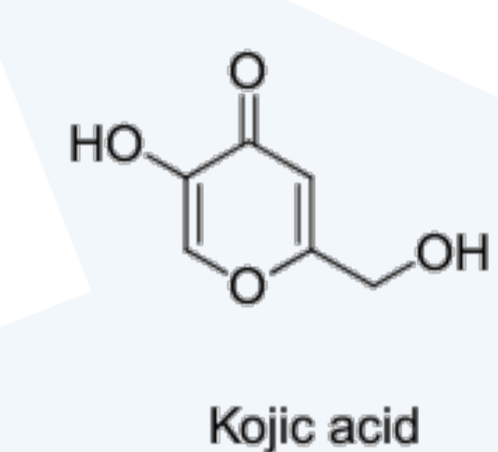
intermediate is O-prenylated at C2 to yield varicoxanthone A, which in turn is C-prenylated to emericillin (Figure 10). The final known step in prenyl-xanthone biosynthesis gives rise to the stereoisomers shamixanthone and epishamixanthone and is catalyzed by XpiC [98]. Alternatively, Nielsen and co-workers include synthesis of angosins by partially reducing the carboxylic acid to an aldehyde, followed by C-prenylation, yielding angosin H and O-prenylation to give angosin A. Subsequent reduction of the aldehyde to a hydroxyl group, and ring closure by dehydration then gives emericillin and shamixanthones [11].

Figure 10. Suggested biosynthesis of the shamixanthons from emodin. Multiple arrows indicate that the number of enzymatic steps are unknown.



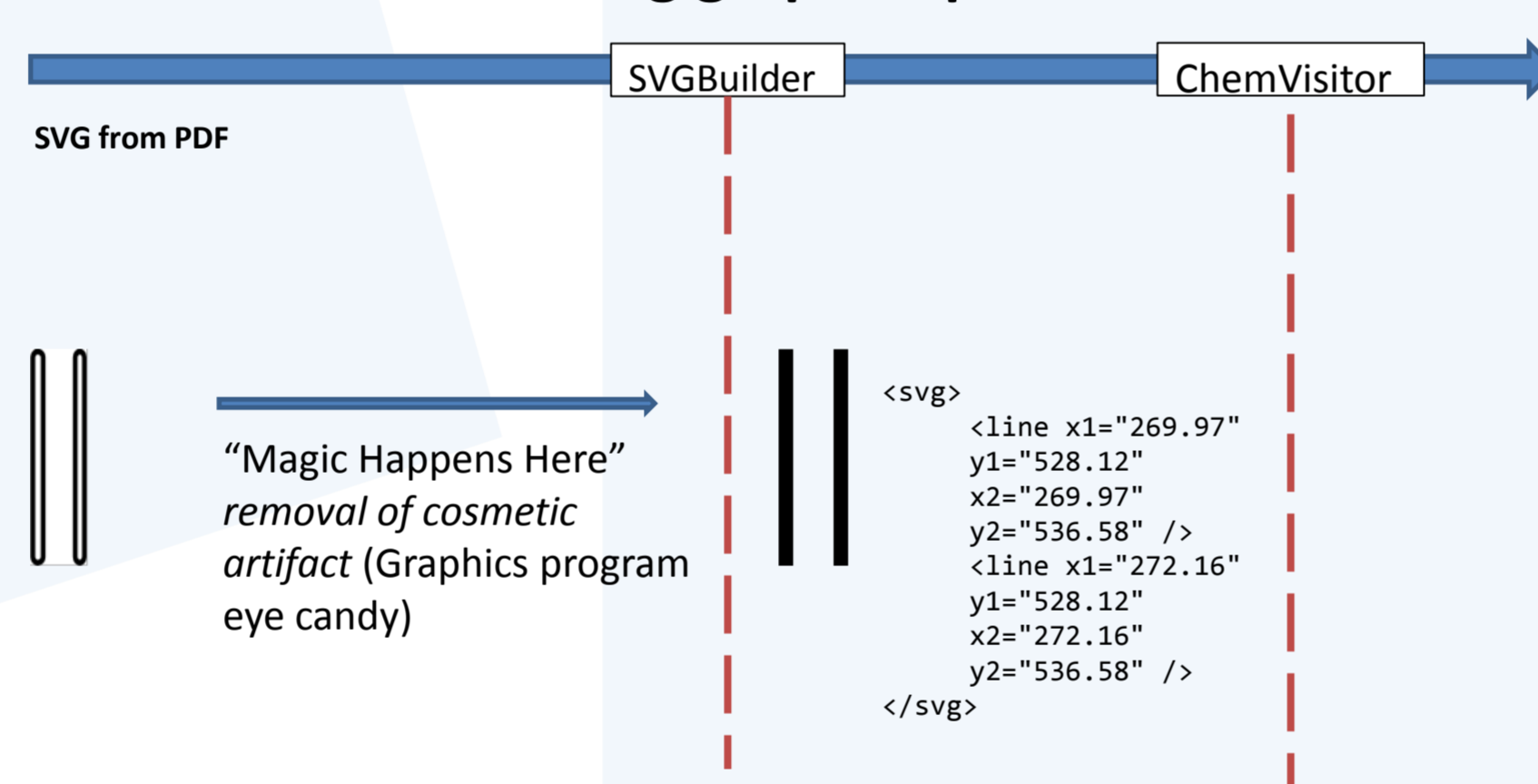
Example diagram from a paper

CML output for one of the molecules in the diagram at the bottom of the poster



```
<?xml version="1.0"?>
<molview xmlns="http://www.xml-cml.org/schema">
  <atomArray>
    <atom id="a1" elementType="C"/>
    <atom id="a2" elementType="C"/>
    <atom id="a3" elementType="H"/>
    <atom id="a4" elementType="H"/>
    <atom id="a5" elementType="H"/>
    <atom id="a6" elementType="H"/>
  </atomArray>
  <bondArray>
    <bond atomRefs2="a1 a2" order="2"/>
    <bond atomRefs2="a1 a3" order="1"/>
    <bond atomRefs2="a1 a4" order="1"/>
    <bond atomRefs2="a2 a5" order="1"/>
    <bond atomRefs2="a2 a6" order="1"/>
  </bondArray>
</molview>
```

Converting graphics primitives



SVG to CML

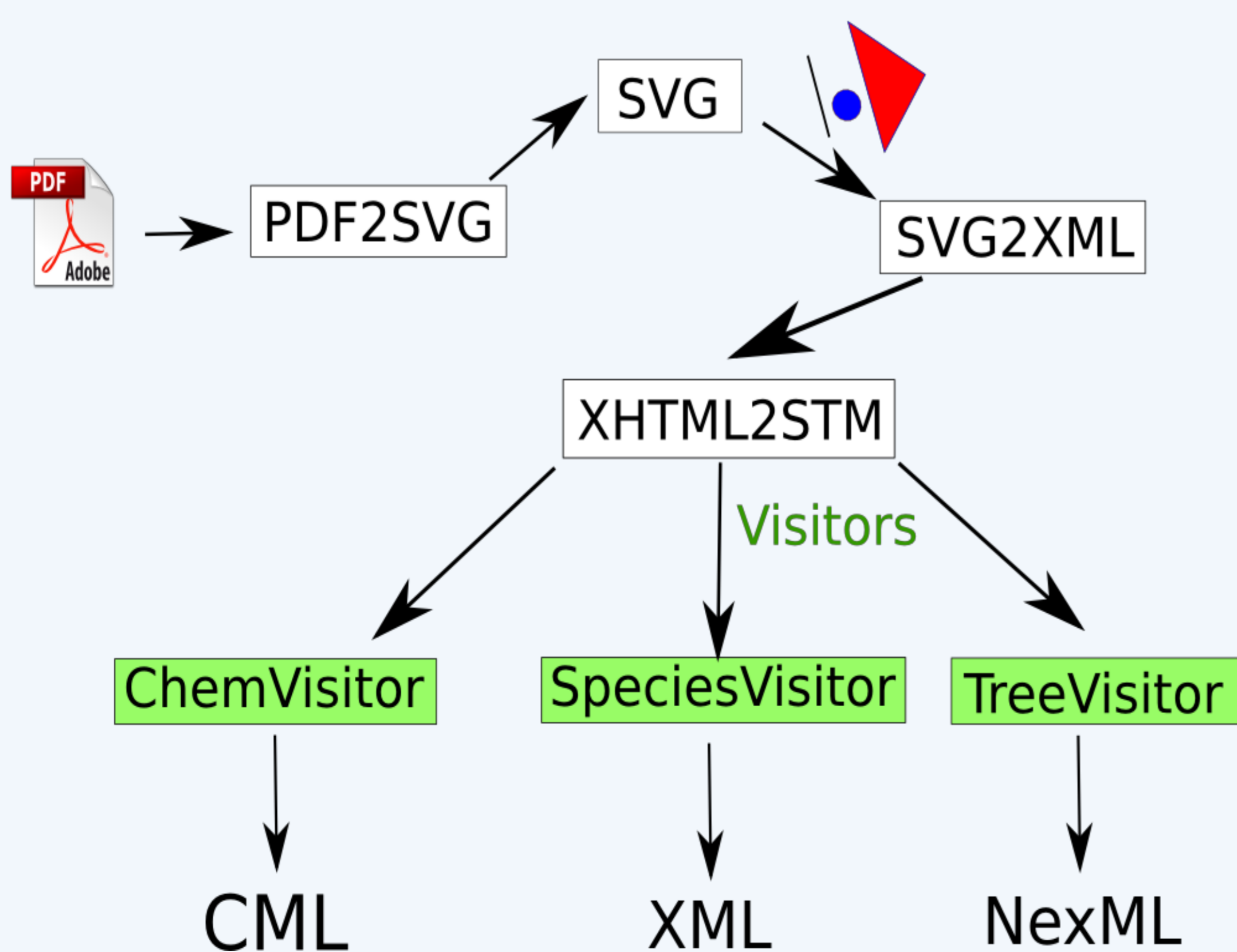
```
<?xml version="1.0"?>
<molview xmlns="http://www.xml-cml.org/schema">
  <atomArray>
    <atom id="a1" elementType="C"/>
    <atom id="a2" elementType="C"/>
    <atom id="a3" elementType="H"/>
    <atom id="a4" elementType="H"/>
    <atom id="a5" elementType="H"/>
    <atom id="a6" elementType="H"/>
  </atomArray>
  <bondArray>
    <bond atomRefs2="a1 a2" order="2"/>
    <bond atomRefs2="a1 a3" order="1"/>
    <bond atomRefs2="a1 a4" order="1"/>
    <bond atomRefs2="a2 a5" order="1"/>
    <bond atomRefs2="a2 a6" order="1"/>
  </bondArray>
</molview>
```

Implicit Hydrogens

Initial processing (currently the same for all visitors)

Chemistry is then determined by identifying which objects are joined to which, and what they are joined with

AMI2 architecture



Other visitors can e.g. extract species names from text and phylogenetic trees from diagrams

The program is available at

<https://bitbucket.org/AndyHowlett/ami2-poc>

Examples

Use of Chemistry Visitor to extract vector diagrams of molecules as CML from a paper

```
ami2poc -i example -o example
ami2poc -i example -o example --target-dir /target-dir
```

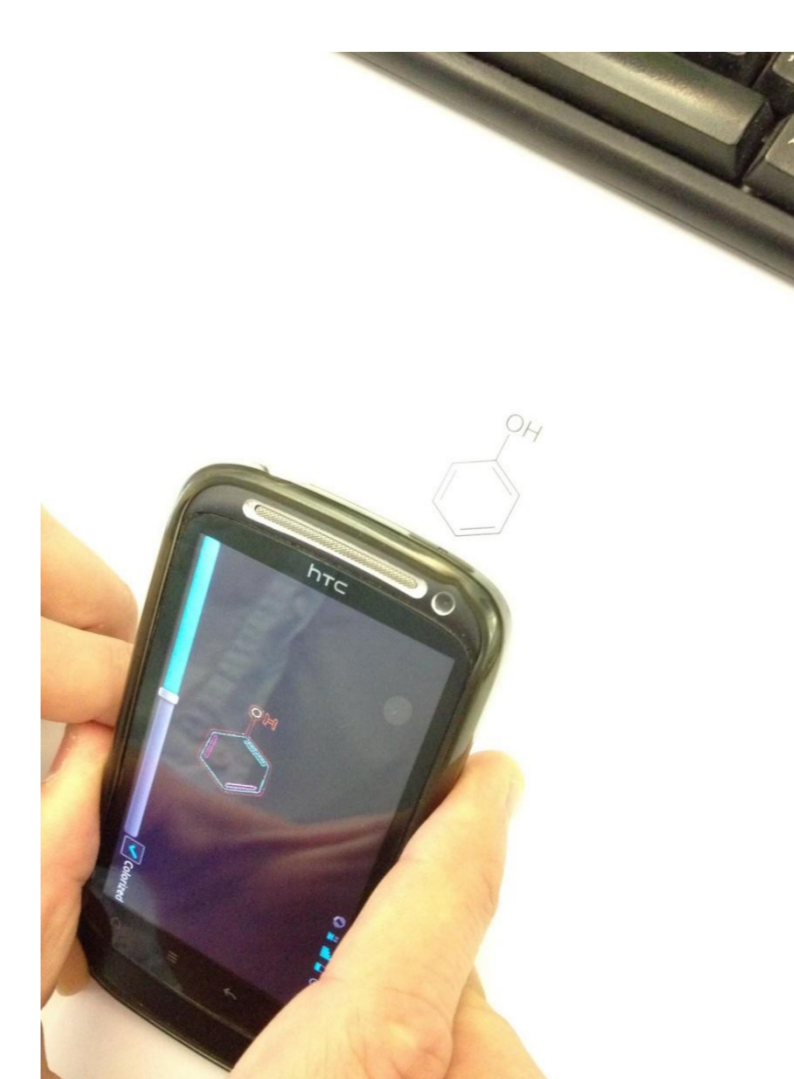
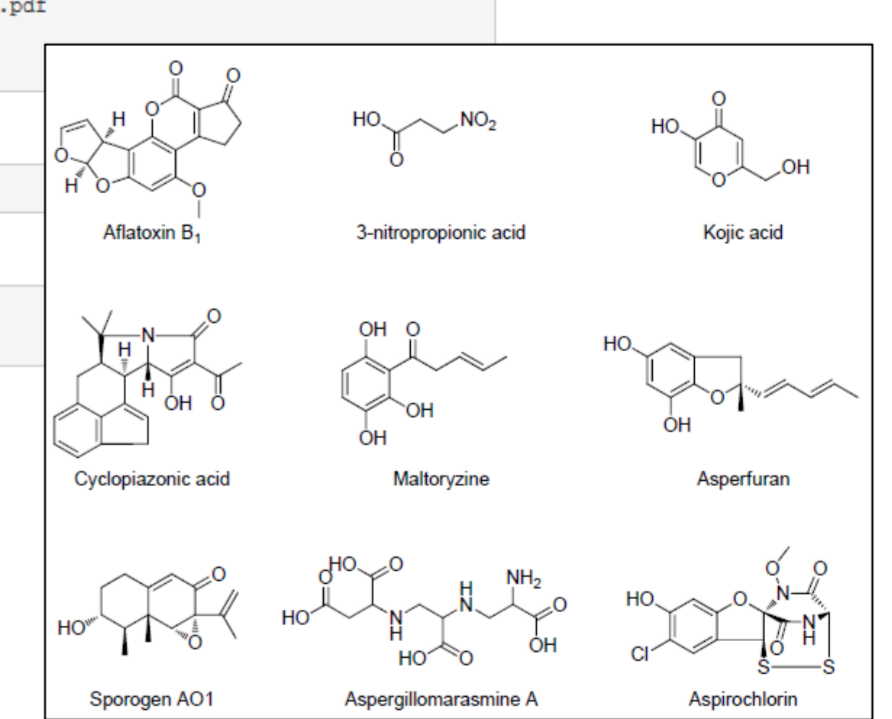
Observe CML output in /target directory

```
ami2poc -i example -o example --target-dir /target-dir
```

View all molecules in figure 1

```
ami2poc -i example -o example --target-dir /target-dir
```

```
ami2poc -i example -o example --target-dir /target-dir
```



We are currently working on applying the same principles to images (either digitally born or from photographs)