

FoLiA: Format for Linguistic Annotation

Maarten van Gompel & Ko van der Sloot

Centre for Language and Speech Technology - Radboud University Nijmegen

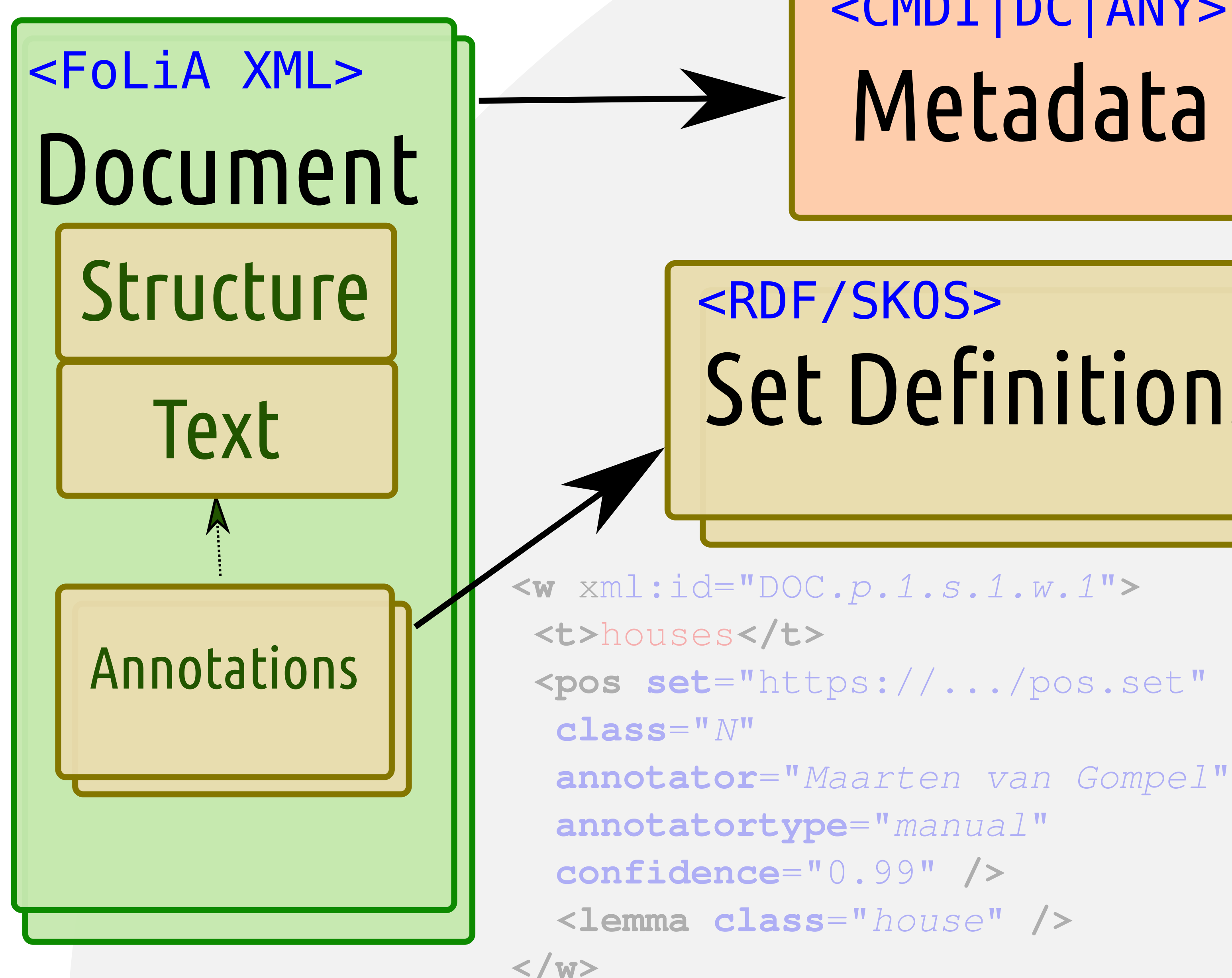


1) Features

<https://proycon.github.io/foia>

- Rich, unified, formalised, XML-based document format for the representation of linguistically annotated resources (incl. corpora)
- Facilitates resource **exchange & interoperability**
- **Specific** support for many linguistic annotation types
- Language & vocabulary (tagset) agnostic: external **set definitions**
- Focus on **practical** usability: lots of open-source **tools** and libraries
- In active use & development for over 7 years

2) Architecture



Structure annotation

Document text/speech structure & tokenisation

Inline annotation & stand-off (span) annotation

various predefined types of linguistic annotation

Higher-order annotation

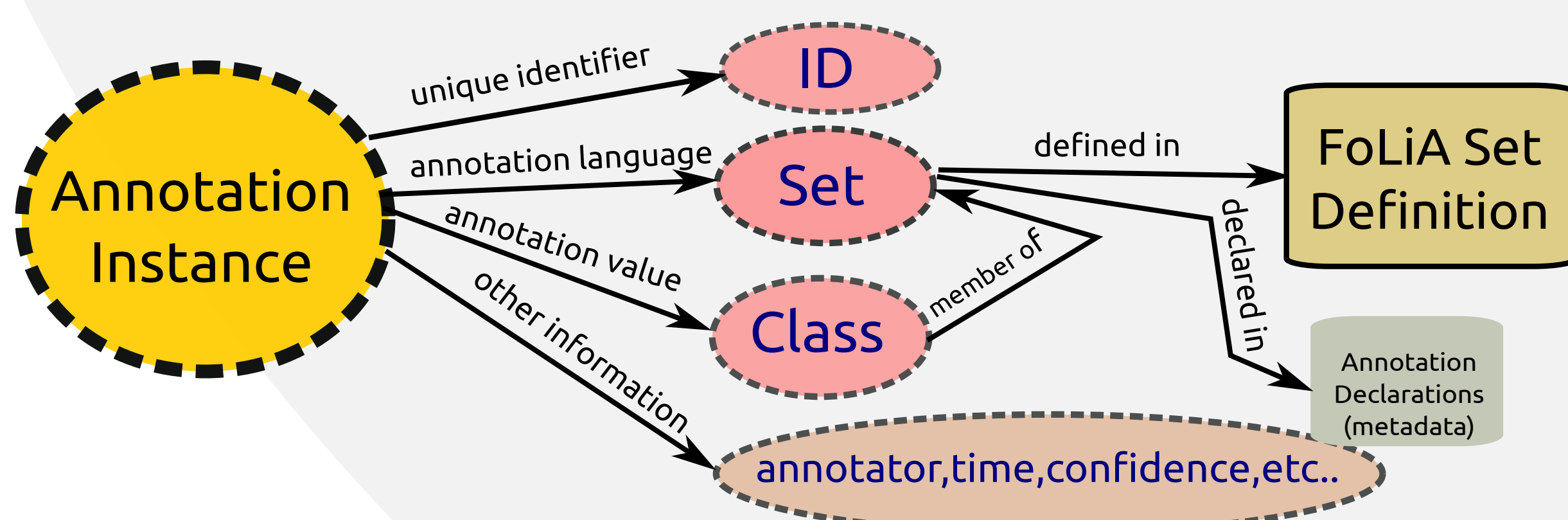
Annotations on annotations

Links to external resources (incl. linked open data)

Generic Attributes

ID, set/class and rich annotator & time information;

```
<pos class="noun">
<feat subset="gender"
class="masc" />
<feat subset="number"
class="sg" />
</pos>
```



3) Tools

- Validators
- Converters
- Visualisation
- Analysis/statistics
- Libraries (Python, C++)

4) Software & Data

- NLP tools (Frog, ucto, TICCL, etc.)
- search tools (MTAS, Black-/Whitelab)
- annotation editors (FLAT)
- corpora (SoNaR, Nederlab, Basilex, etc.)

5) Python Library

```
$ pip install pynlpl
> from pynlpl import folia
> doc = folia.Document(
    file="doc.xml")
> for word in doc.words():
    . . .
```

FoLiA, its main libraries and tools are all licensed under the GNU Public Licence v3 - Funding by CLARIN-NL & CLARIAH

