# 10. Text Annotation Structure

# Text Annotation Structure

- When applying any linguistic resource to a text, NooJ adds, or removes annotations to the Text Annotation Structure (TAS)

- During the lexical analysis, NooJ applies dictionaries and morphological grammars to the text, to add annotations that represent ALUs into the TAS

- During the syntactic analysis, NooJ applies syntactic grammars to the text to add structural annotations, or remove ALU annotations (*e.g.*, when solving ambiguities).

# Text Annotation Structure

NooJ lexical analyzer can annotate any ALU and represent them in the TAS:

- Agglutinated or contracted forms:

    *... cannot ...* ⟶ **<can,V> <not,ADV>**

- Simple words:

    *... is ...* ⟶ **<be,V+PR+3+s>**

- Multiword units:

    *... Blue collars ...* ⟶ **<blue collar,N+Hum+p>**

- Discontinuous expressions:

    *turns* the light *off* ⟶ **<turn off,V+PR+3+s>**

# Text Annotation Structure

NooJ uses linguistic resources to recognize ALUs:

- Morphological grammars:

cannot/<can,V><not,ADV>

*... cannot ...* → **<can,V> <not,ADV>**

- Dictionaries associated with inflectional/derivational grammars:

`be,V+FLX=BE+Aux`

`blue collar,N+Hum+FLX=TABLE`

*... is ...* → **<be,V+Aux+PR+3+s>**

*... Blue collars ...* → **<blue collar,N+Hum+p>**

- Dictionary/grammar pairs (see later)

`turn,V+PV+FLX=HELP+Part="off"`

<V+PV> <WF>* $PV$Part

*... turns* the light *off ...* → **<turn off,V+PR+3+s>**

# Annotation of the wordform *cannot*

# Annotation of the multiword unit *all of a sudden*

# Annotation of the discontinuous expression *take … back*

# Automatic Lexical Analysis
# Representing lexical and morphological ambiguities

# Text Annotation Structure

## The TAS can be exported as an XML file

# Text Annotation Structure

The TAS can be exported as an XML file, but…

- Consider the two multiword units, represented in NooJ dictionary:

`all of a sudden,ADV`

`sudden death syndrome,NOUN`

*… all of a sudden death syndrome …*

- When applying this dictionary to the above text, NooJ will represent the ambiguity in this TAS, as any of these two multiword units might be occurring.

- But XML cannot handle crossed-scoped annotations such as:

*… <ADV>all of a <NOUN>sudden</ADV> death syndrome</NOUN> …*

- **Solution**: Either remove all ambiguities before exporting the TAS, or give priority to the first occurring multiword unit, e.g.,

*… <ADV>all of a sudden</ADV> death syndrome …*

# Text Annotation Structure

## Removing ambiguities

- Removing ambiguities = deleting annotations from the TAS

- By exploring contexts. For example "can" is ambiguous, but not in these two contexts:

    *… They can* (VERB) *open the can* (NOUN) *of beer …*

- We will use syntactic or semantic grammars that describe contexts in which some ambiguities can be solved.


- However, in general, many ambiguities will not be solved, e.g.,

    *… There is a round table in room A32 …*

    (a meeting, or a round piece of furniture?)

# CONGRATULATIONS

You know how to automatically annotate a text by applying dictionaries and morphological grammars to the text.