# 14. Text Annotation Structure (Advanced)

# 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

- When applying syntactic grammars to a text, NooJ can add annotations (e.g., structural), or remove annotations (*e.g.*, lexical hypotheses).

# Syntactic Analysis
## Text, right after the lexical analysis

*man*: noun or verb (several forms); *eat*: noun or verb

# Syntactic Analysis
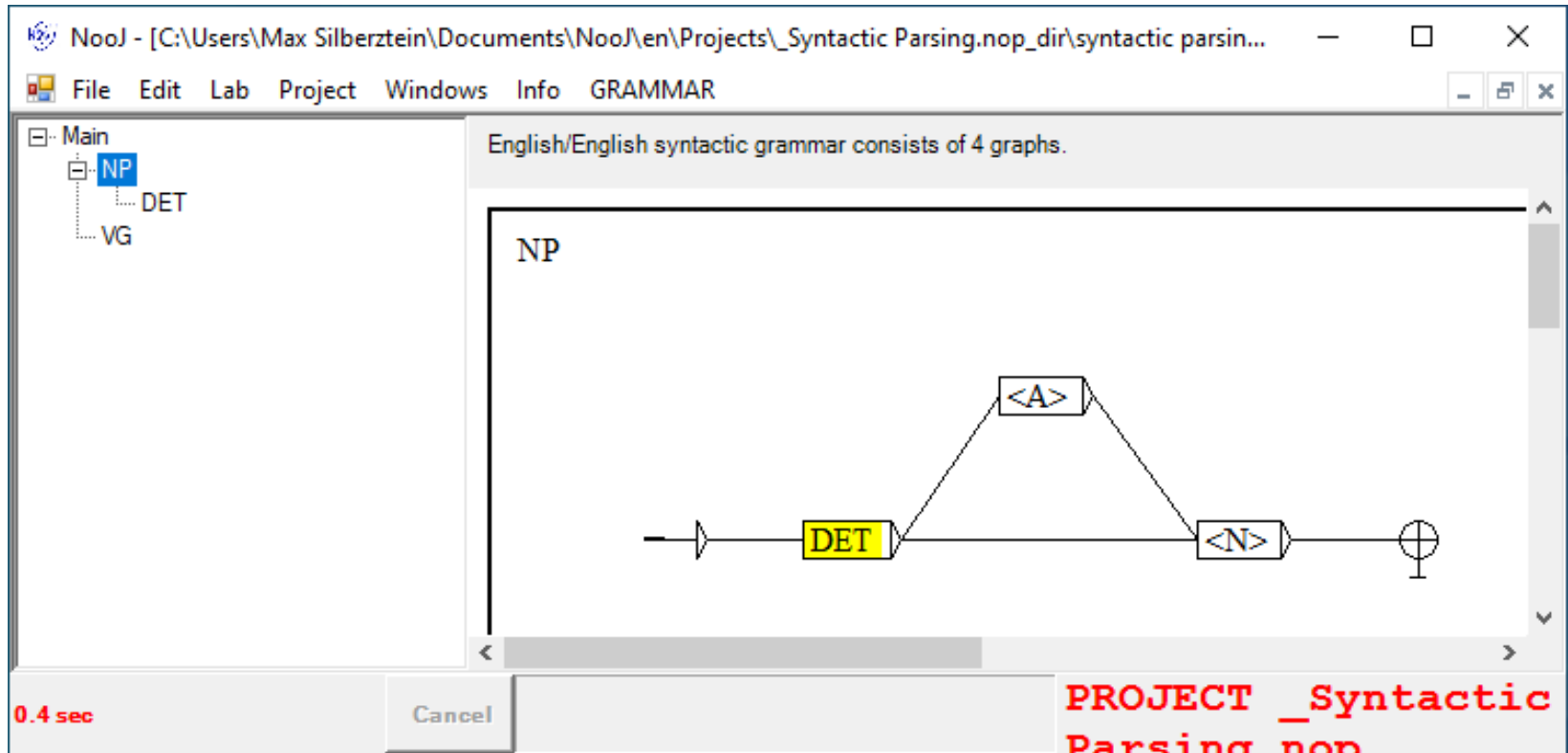## Syntactic grammar produces structural annotations



- Annotation <SENTENCE>
- Annotation <NB+Sub>
- Annotation <VG>
- Annotation <NB+Obj>

# Syntactic Analysis
## Syntactic grammar produces structural annotations



- ## Several embedded grammars

# Syntactic Analysis
## Applying grammar to text

# Syntactic Analysis
## Applying grammar to text



The output represents the structure produced by the annotations

# Syntactic Analysis

## CONCORDANCE > Display Syntactic Analysis

NooJ displays the parse tree:

# Syntactic Analysis
## CONCORDANCE > Display Syntactic Analysis

The TAS is displayed as a tree:

# Syntactic Analysis
## Parse *vs.* Structural trees

- The **parse tree** represents the structure of the grammar, rather than the structure of the sentence.

➤ It is useful to debug a grammar, as it shows how the grammar was explored during parsing

- The **structural tree** represents the TAS. It is produced by annotations in the grammar, and is independent from the structure of the grammar.

➤ It is useful to accumulate and share grammars, as it is independent from how linguists want to organize their grammar; it also allows NooJ to optimize grammars (*e.g.*, remove useless rules and recursions, etc.) without any consequence

# Syntactic Analysis
## Atomic Linguistic Units

- There are four types of ALUs: affixes, simple words, multiword units and discontinuous expressions.

➢ ALUs are represented by annotations in the TAS

- Syntactic trees must represent the ALUs

# Syntactic Analysis
# Contracted and agglutinated forms

*In French, aux* is a contracted form of *à les*

# Syntactic Analysis
# Contracted and agglutinated forms

*In French, aux* is a contracted form of *à les*

# Syntactic Analysis
## Discontinuous expressions

- *Il ne baisse pas le ton* [he does not lower his voice]
- In French, the negation *ne ... pas* is discontinuous.
- In French, the frozen expression *baisser ... le ton* is discontinuous.

**The TAS:**



```
Language is "French (fr)".
Text Delimiter is: \n (NEWLINE)
Text contains 2 Text Units (TUs).
7 tokens including:
6 word forms
Text contains 46 annotations (74 different)
```

nom | nepas,NEG | baisser,V+C1D+Temps=PR+N0=N0Hum+DET1=le+N1=ton+Pers=3+Nb=s | ADV | DET | N+C1D

# Syntactic Analysis
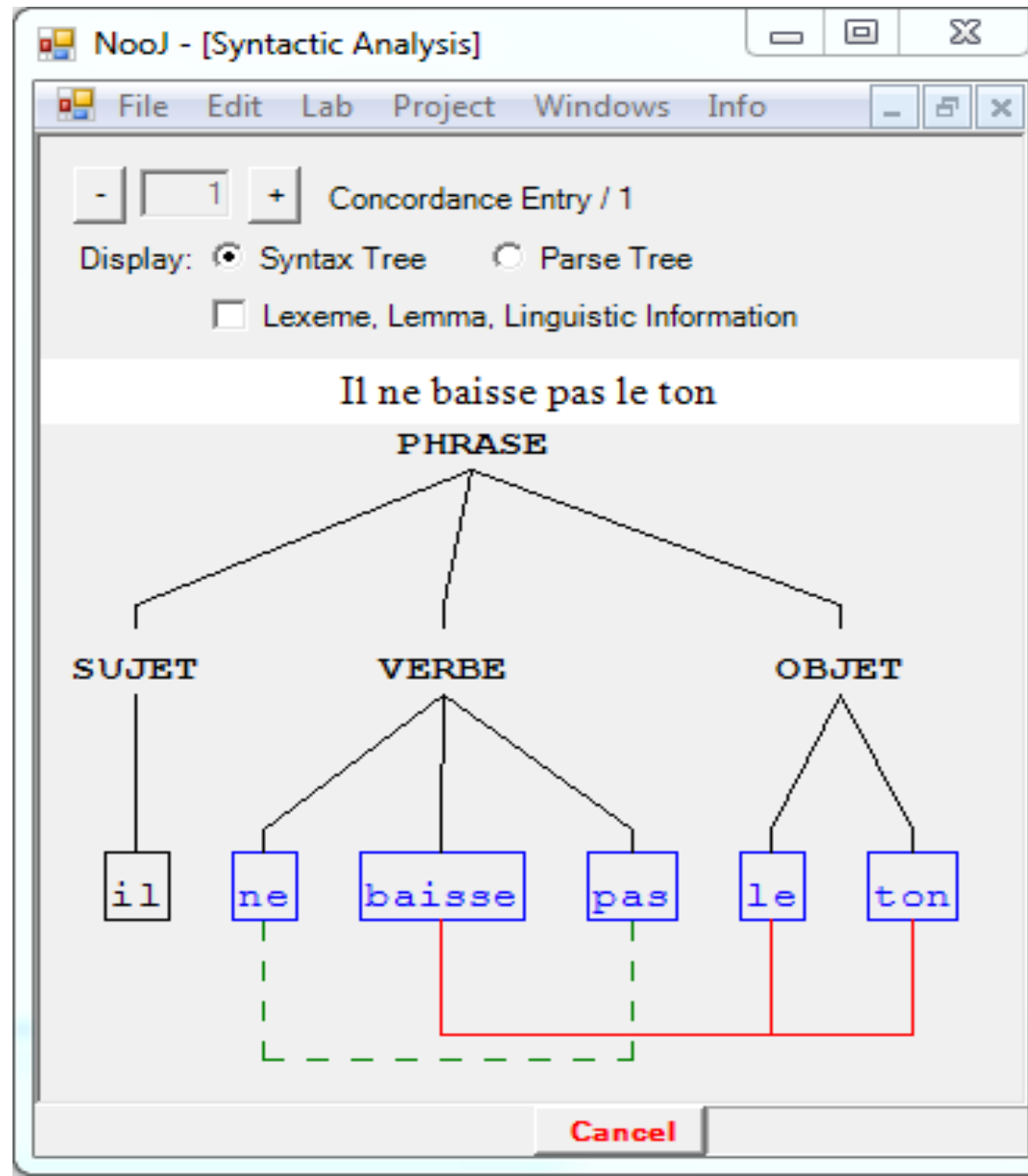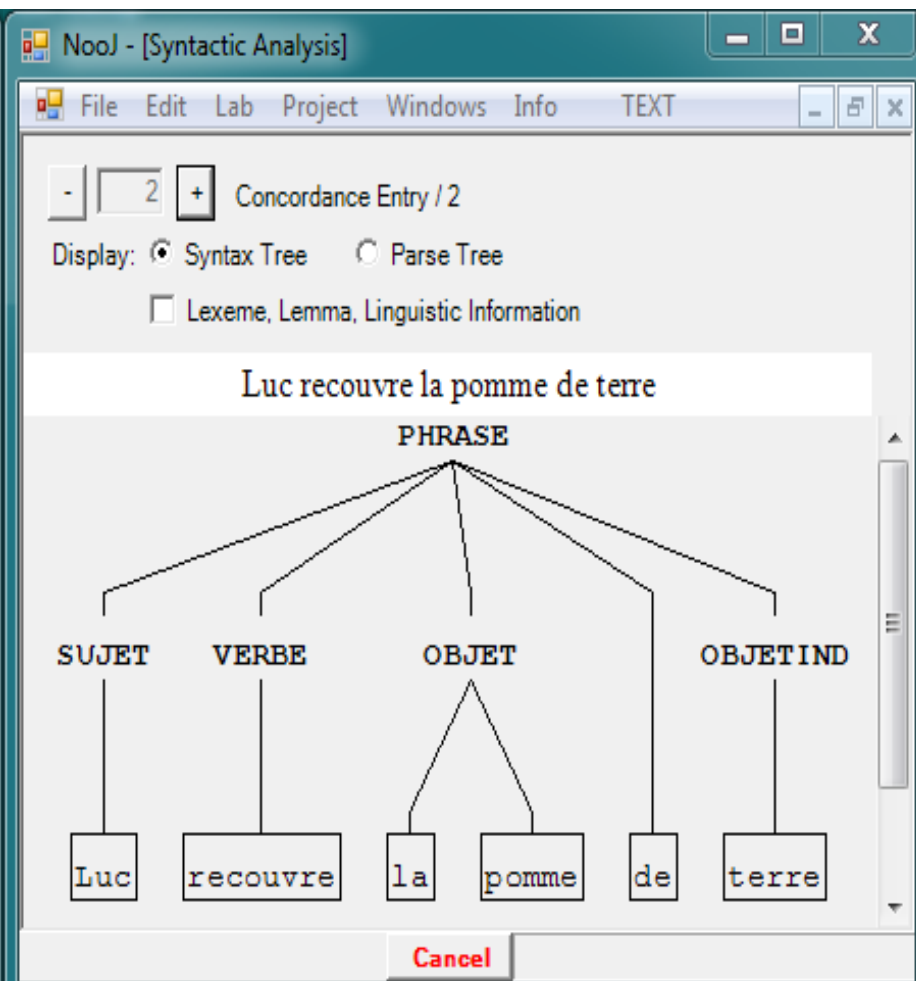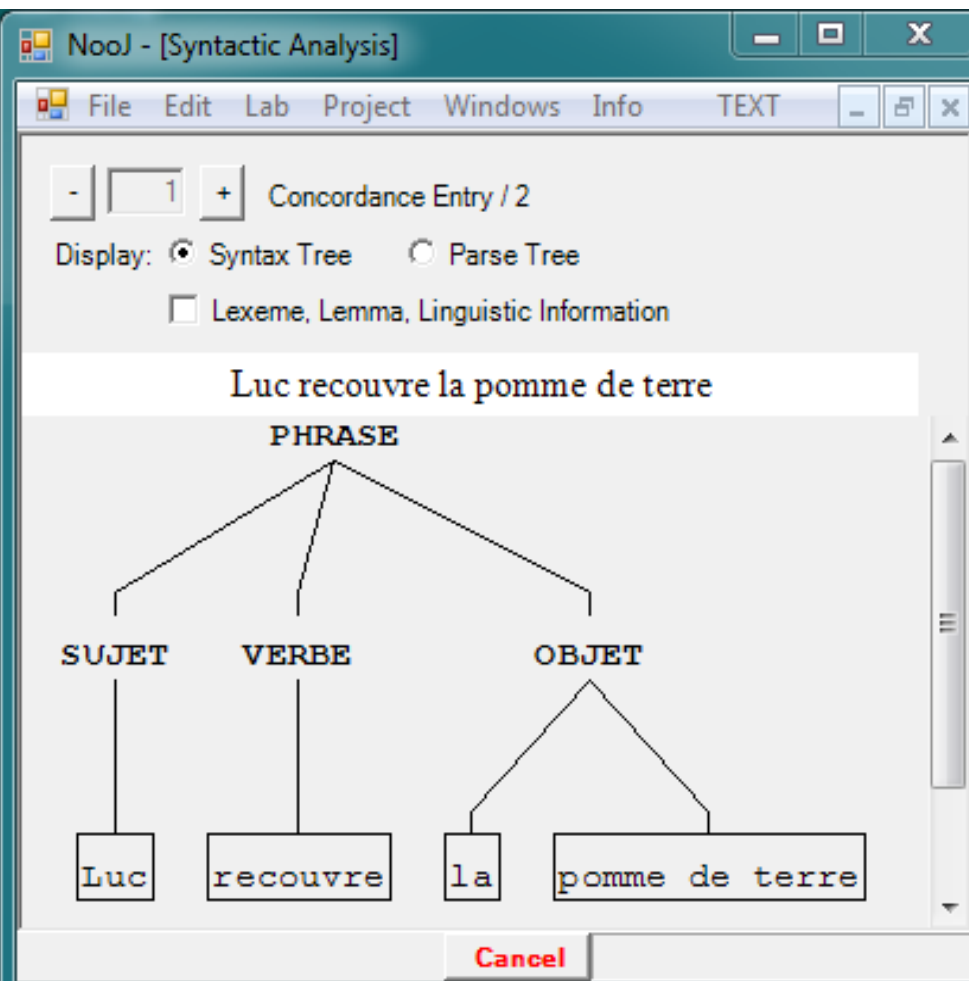# Discontinuous expressions

- The negation *ne ... pas* is discontinuous.

- The frozen expression *baisser ... le ton* is discontinuous.

# Syntactic Analysis
# Ambiguities

*Luc recouvre la pomme de terre*
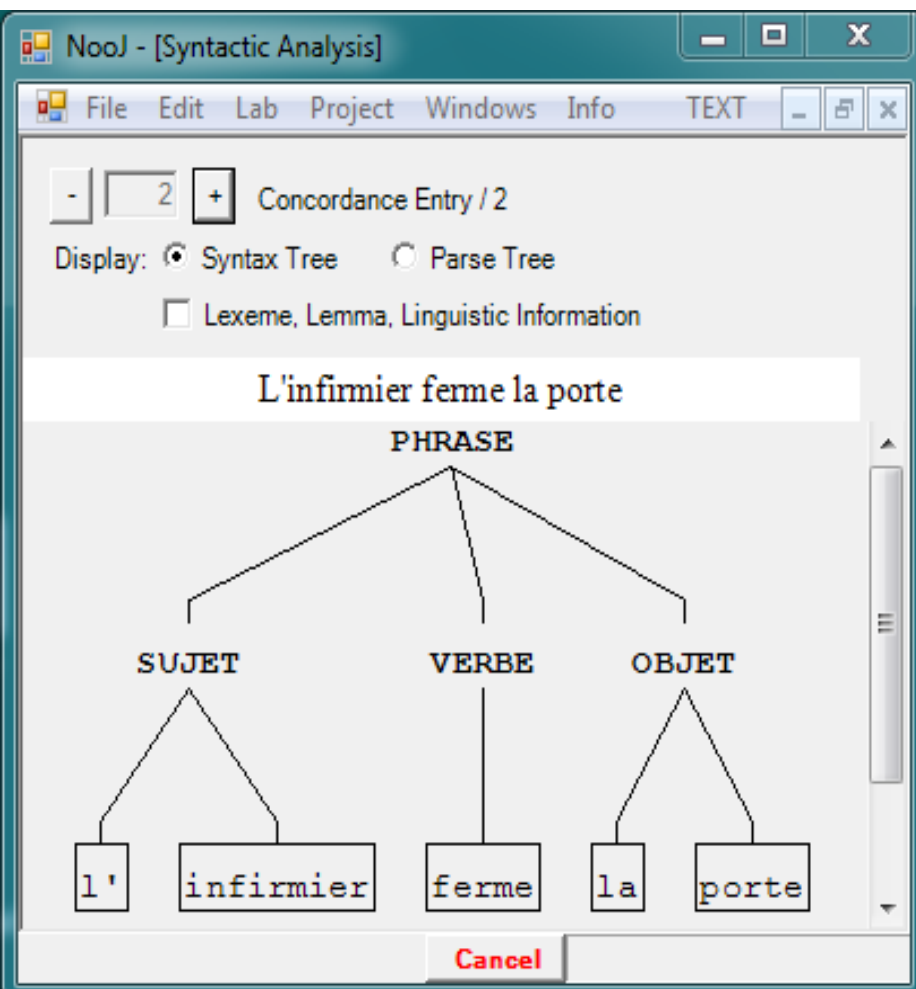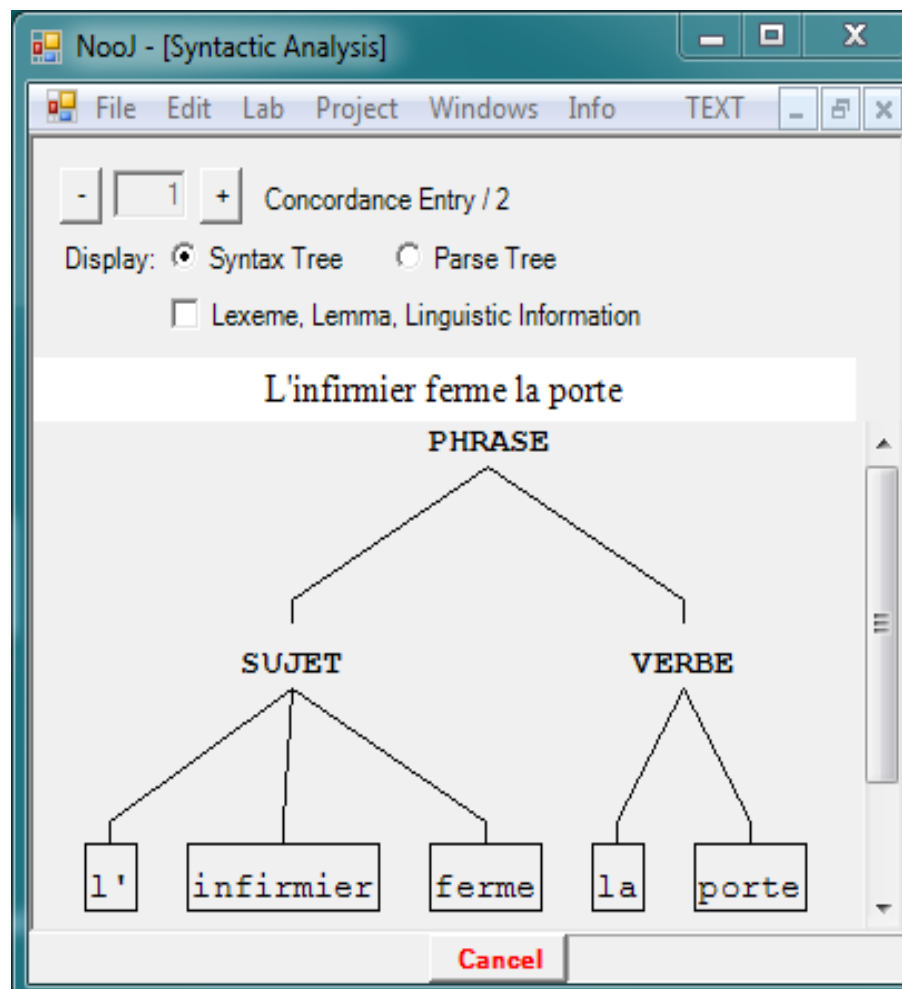*[Luc covers the apple with earth]* or *[Luc covers the potato]*
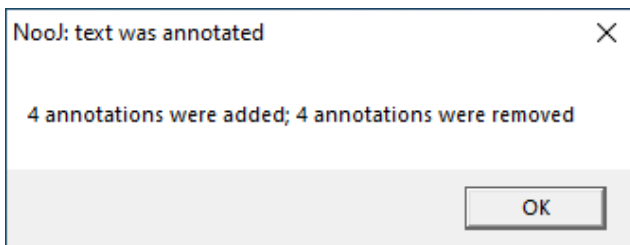
# Syntactic Analysis
# Ambiguities

*L'infirmier ferme la porte*
*[The nurse closes the door]* or *[the firm nurse carries her]*

# Syntactic Analysis
# Disambiguation

CONCORDANCE > Annotate Text (add/remove annotations)

# Syntactic Analysis
# Disambiguation: before and after

# CONGRATULATIONS

You know how to perform various lexical, morphological, syntactic and semantic analyses by annotating texts with various types of information