

Software and the digital humanities (N. Derschowitz & W. Sack)  
Institut d'Études Avancées Paris 4 novembre 2015



# **GRAIL**

## **A CATEGORIAL PARSER FOR SYNTAX AND SEMANTICS**

**Discourse on software vs software on discourse**

**A logician's approach: discourse on software on discourse**

**CHRISTIAN RETORÉ UNIVERSITÉ DE MONTPELLIER & LIRMM**  
**(WITH RICHARD MOOT CNRS LABRI & LIRMM)**

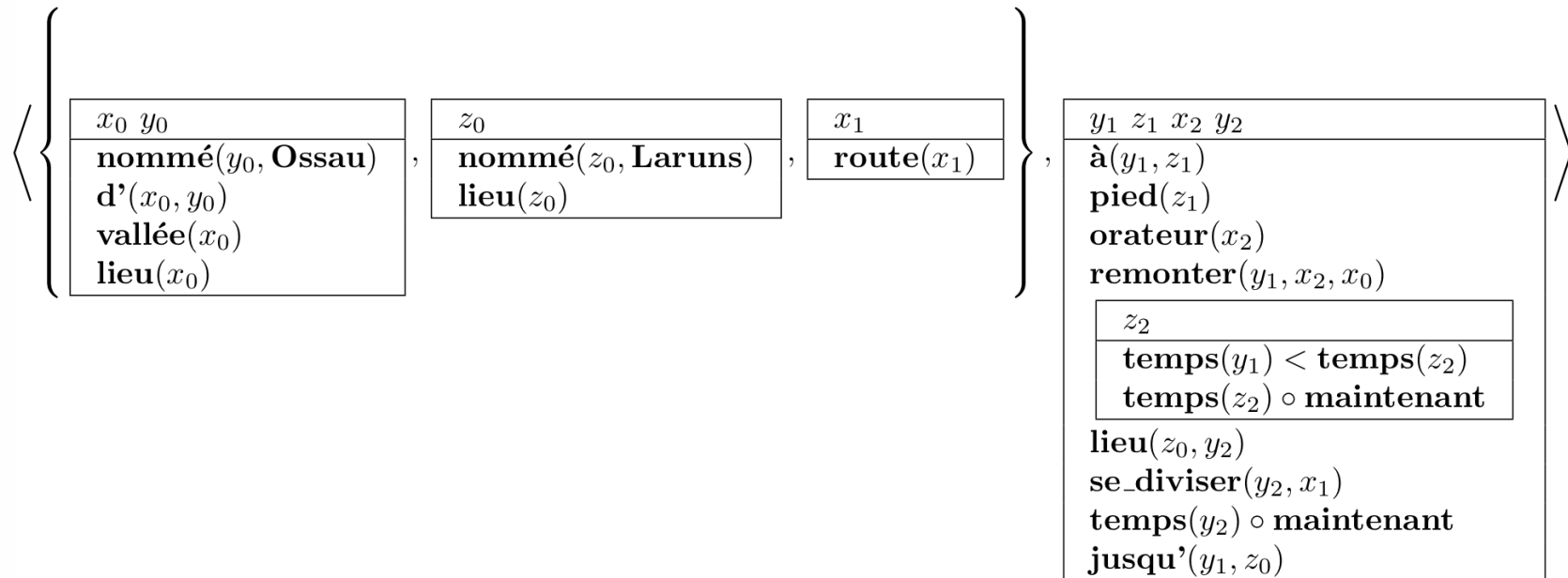


# FROM TEXT TO LOGIC

- An example:
  - Children will have a pizza.
  - For all child x there exists a pizza ordered for x
  - There exists a pizza p ordered for all children.
- **what is asserted**  $\neq$  what it speaks about
  - Ok it's windy but that's not an ouragan.
  - Was Geach a student of Wittgenstein?
    - En 1941, IL épousa la philosophe Elizabeth Anscombe, grâce à LAQUELLE IL entra en contact avec Ludwig Wittgenstein. Bien qu'IL N'ait JAMAIS suivi l'enseignement académique de CE DERNIER, cependant IL EN éprouva fortement l'influence.



# WHAT DOES SEMANTICS LOOKS LIKE



# SYNTACTIC ANALYSIS

- Categorical Grammar
- Rules are language independent
- The lexicon makes the difference
- Word  $\rightarrow$  categories  $\sim$  logical formulae
- Parsing as deduction





# SYNTAX ➤ SEMANTICS

- Grammatical categories ~ logical structure
  - Common noun : property
  - Verb phrase: property
  - Definite noun phrase: individual



# COMPOSITIONAL SEMANTICS

- Frege:  
the meaning of a compound expression is a function of the meanings of its parts
- Montague: **AND** of the syntactic structure
- Semantic lambda terms in the lexicon combined according to syntax ...
- reduction (computation)  
-> logical formula (~meaning?)



# LEXICAL SEMANTICS WORD MEANING IN CONTEXT

- Word meaning in context (polysemy)
  - The bank is muddy.
  - The bank resisted the crisis.
- Types & sorts to force coherence:
  - “bark” applies to animals (usually dogs)
  - \* The chair barked
  - The secretary barked upon students that registered lately.
- Copredication
  - \* On Sunday, Barcelona won and voted.
  - Barcelona voted and will host the Olympiads.





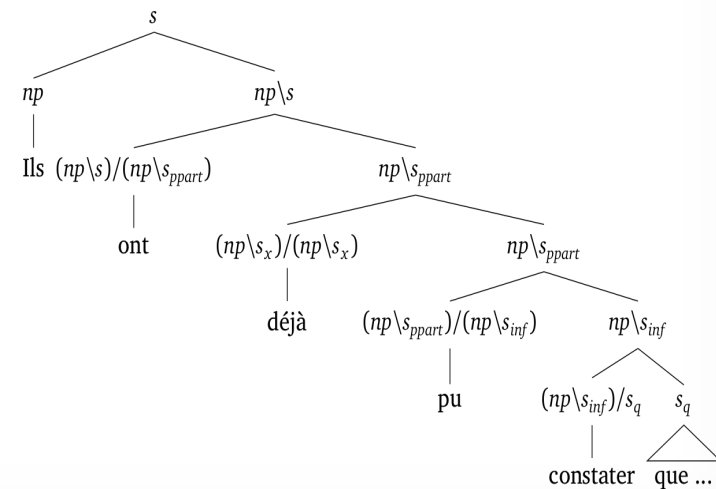
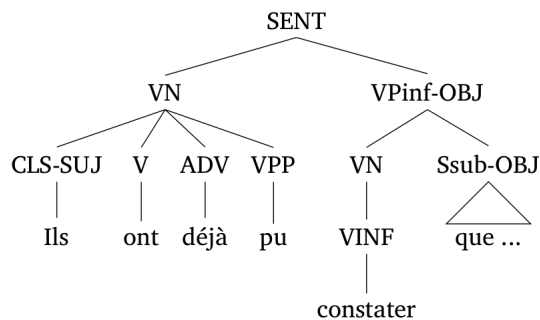
## OTHER ASPECTS OF LINGUISTIC CONTEXT

- Reference and anaphora
  - Anaphora resolution (not implemented)
  - Carlos' dog thinks he does not like him.  
he≠him
  - If a farmer owns a donkey then he beats it.  
(accessibility is implemented)
- Presupposition (implemented)
  - A thinks/knows B came. => B came? no/yes



# ACQUISITION OF THE LINGUISTIC RESOURCES

- Syntax (ok)
- Compositional semantics (little)
- Lexical semantics (+/-)
- Discursive structures (little)



# GRAIL FOR FRENCH

- Initially developed for spoken Dutch
- Acquired from annotated corpora  
constituents and dependencies  
corpus Paris 7 -> tигра annotations --> categories
- Up to 80 categories per word
- Supertagging
- Parsing



# A HISTORICAL & REGIONAL CORPUS

- Travel stories in the Pyrenees
- Set of XIX century texts
- 576 334 words
  
- A difficult question: can we reconstruct the itineraries followed by the narrator?



# CASE STUDY: FICTIVE MOTION AND GRAIL

- cette route qui monte sans cesse
- For all x if x follows the road then x keeps going up

$$\begin{array}{c}
 \frac{\text{monte} \quad [Lex] \quad \frac{\text{sans} \quad [Lex] \quad \frac{\text{cesse} \quad [Lex]}{n} \quad [E]}{(s \setminus_1 s)/n} \quad [E]}{np \setminus s_{main}} \quad [E]}{[q_0 \vdash np]^1} \quad [E]}{q_0 \circ \text{monte} \vdash s_{main}} \quad [E]} \\
 \frac{[p_0 \vdash np]^0}{p_0 \circ (\text{monte} \circ_1 (\text{sans} \circ \text{cesse})) \vdash s} \quad [\wedge]_1}{\text{monte} \circ_1 (\text{sans} \circ \text{cesse}) \vdash np \setminus s_{main}} \quad [E]} \\
 \frac{p_0 \circ (\text{monte} \circ_1 (\text{sans} \circ \text{cesse})) \vdash s_{main}}{p_0 \circ ((\text{monte} \circ_1 (\text{sans} \circ \text{cesse})) \circ_1 (\text{pendant} \circ (\text{deux} \circ \text{lieues}))) \vdash s} \quad [\wedge]_0}{(\text{monte} \circ_1 (\text{sans} \circ \text{cesse})) \circ_1 (\text{pendant} \circ (\text{deux} \circ \text{lieues})) \vdash np \setminus s_{main}} \quad [E]} \\
 \frac{\text{qui} \quad [Lex]}{(n \setminus n)/(np \setminus s_{main})} \quad [E]}{\text{route} \quad [Lex]} \quad \frac{n}{n} \\
 \frac{\text{cette} \quad [Lex]}{np/n} \quad [E]}{\text{route} \circ (\text{qui} \circ ((\text{monte} \circ_1 (\text{sans} \circ \text{cesse})) \circ_1 (\text{pendant} \circ (\text{deux} \circ \text{lieues})))) \vdash n} \quad [E]} \\
 \frac{\text{sur} \quad [Lex]}{(s \setminus_1 s)/np} \quad [E]}{\text{cette} \circ (\text{route} \circ (\text{qui} \circ ((\text{monte} \circ_1 (\text{sans} \circ \text{cesse})) \circ_1 (\text{pendant} \circ (\text{deux} \circ \text{lieues})))) \vdash np} \quad [E]} \\
 \frac{\text{avance} \quad [Lex]}{np \setminus s_{main}} \quad [E]}{\text{sur} \circ (\text{cette} \circ (\text{route} \circ (\text{qui} \circ ((\text{monte} \circ_1 (\text{sans} \circ \text{cesse})) \circ_1 (\text{pendant} \circ (\text{deux} \circ \text{lieues})))) \vdash np} \quad [E]} \\
 \frac{\text{On} \quad [Lex]}{np} \quad [E]}{\text{avance} \circ_1 (\text{sur} \circ (\text{cette} \circ (\text{route} \circ (\text{qui} \circ ((\text{monte} \circ_1 (\text{sans} \circ \text{cesse})) \circ_1 (\text{pendant} \circ (\text{deux} \circ \text{lieues})))) \vdash np \setminus s_{main}} \quad [Wpop_{vp}]} \\
 \frac{\text{On} \circ (\text{avance} \circ_1 (\text{sur} \circ (\text{cette} \circ (\text{route} \circ (\text{qui} \circ ((\text{monte} \circ_1 (\text{sans} \circ \text{cesse})) \circ_1 (\text{pendant} \circ (\text{deux} \circ \text{lieues})))) \vdash np \setminus s_{main}} \quad [E]} \\
 4. \quad \text{On} \circ (\text{avance} \circ_1 (\text{sur} \circ (\text{cette} \circ (\text{route} \circ (\text{qui} \circ ((\text{monte} \circ_1 (\text{sans} \circ \text{cesse})) \circ_1 (\text{pendant} \circ (\text{deux} \circ \text{lieues})))) \vdash s_{main}}
 \end{array}$$



# GRAIL: NUMBERS

- Swi Prolog (interf. Tcl/Tk) acquisition+parsing 30.000 lines
- Supertagging: C software by Clark & Curran
- Corpus French Treebank 423 152 words
- Lexicon
  - 28753 distinct words
  - Up to 80 formulae per word (et, être, prepositions...)
- Semantic lexicon:
  - 350 default rules
  - 500 specific entries (coordination, quantifiers)



# APPLICATION OF LOGICAL ANALYSIS

- Analysis of argumentative dialog
- Text entailment
  - One paragraph  $\rightarrow$  logical formulae  $\Gamma$
  - One sentence  $\rightarrow$  a logical formula  $C$
  - Is  $C$  a consequence of  $\Gamma$  ?
  - (Coq Coquand, **Huet**, **Dowek**,...)



# RELEVANCE TO THE WORKSHOP : TRANSLATION

- Translation of text into logic (faithfulness, approximation)
  - CS translation :  
annotated corpora -> CG trees (data)
  - Sentence +supertagging -> parse tree
  - Parse tree + lexicon -> semantics





## RELEVANCE TO THE WORKSHOP : TRANSLATION

- Le sens ne se produit jamais que de la traduction d'un discours en un autre. (Lacan)
- Translation of linguistic/cognitive theories into principles of computational analysis (artificial intelligence ?)
- Approximation:
  - Degrees of grammaticality, of assertion, etc.
  - « un peu » vs. « peu » (Anscombe & Ducrot)



# RELEVANCE TO THE WORKSHOP

- Do there exist humanities-based approaches to software that acknowledge that software is sometimes not a tool, but rather a condition or a problem?
- Software implements/test linguistic/cognitive theories (here on syntax and semantics)



# A PARTISAN PARSER

- A view of syntax, semantics, discourse theories:
  - **Language as a structure and as a vehicle for information**
  - vs. Language as opinion expression and speech acts
- Implicit claim (cf. Turing test or Chomsky): language faculty is a **computational** device

