

Alpino: A Wide Coverage Computational Grammar for Dutch

Gosse Bouma and Gertjan van Noord
Alfa-informatica, RUG

CLIN 00

Overview

- PIONIER project 'Algorithms for Linguistic Processing'
- The Alpino grammar,
- Lexical resources,
- Construction of Dependency Trees,
- Treebank and evaluation,
- Future work.

Algorithms for Linguistic Processing

- Efficient processing and disambiguation of natural language,
- Develop wide-coverage Dutch grammar,
- Study disambiguation techniques,
- Evaluate coverage & disambiguation,
- *(Efficiency & Finite-state approximations).*

Grammar

- Lexicalized (HPSG-style) grammar,
- Extension of the NWO-TST (*OVIS*) grammar,
- Added rules for written language,
- Incorporated lexical entries based on Celex and Parole.

Rule Coverage

- Sentence types: declaratives, yes/no & WH-questions, topicalization, imperatives, subordinate clauses,
- NPs: relatives, sbar-complements, titles (*minister zalm*), measure phrases, temporal expressions, ..
- VP syntax: NP, PP, VP, SBAR complements, predicative phrases, adjuncts, verb clusters, particles, passives.
- Coordination of maximal projections (NP, PP, S, ...).

Inheritance in Rule Definitions

- 114 rules
- `pp --> p np is-a head-comps-struct.`
- `head-comps-struct is-a headed-struct.`
- `headed-struct` satisfies
 - ★ `head-feature` *principle*,
 - ★ `valence` *principle*,
 - ★ `filler` *principle*.

Example Rule

$$\text{rule}(\underset{pp}{\left[\begin{array}{cc} \text{wh} & \boxed{1} \\ \text{slash} & \langle \rangle \\ \text{prep} & \boxed{2} \\ \text{ppost} & no \\ \text{dt} & \boxed{3} \end{array} \right]}, \langle \underset{p}{\left[\begin{array}{cc} \text{sc} & \langle \boxed{4} \rangle \\ \text{wh} & \boxed{1} \\ \text{prep} & \boxed{2} \\ \text{ppost} & no \\ \text{dt} & \boxed{3} \end{array} \right]}, \underset{np}{\boxed{4}} \left[\text{nform} \quad \neg er \right] \rangle).$$

Inheritance for Lexical Entries

- *'toerekenen'* **is-a** trans-particle-verb
- trans-particle-verb **is-a** trans-verb
- trans-verb **is-a** np-subj-verb
- np-subj-verb **is-a** verb
- verb **is-a** lexical-sign.
- lexical-sign satisfies argument-realization.

Recursive Constraints & Co-routining

- Slash-introduction defined as a constraint on mapping from DEPENDENTS (and SUBJ) to SUBCAT and SLASH (Bouma, Malouf, Sag, 2001).
- Verb-raising verbs defined using argument-inheritance (append of SUBCAT-lists) (Bouma and van Noord, 97),
- Co-routining is used for implementation of such constraints (van Noord and Bouma ,1994).

$$\text{lex}(\text{v} \left[\begin{array}{ll} \text{deps} & \textcircled{1} \left\langle \begin{array}{ll} \text{case} & \text{acc} \\ \text{nform} & \text{norm} \end{array} \right\rangle_{np} \\ \text{subj} & \textcircled{2} \left[\begin{array}{ll} \text{agr} & \text{sg\&thi} \\ \text{case} & \text{nom} \\ \text{nform} & \text{norm} \end{array} \right]_{np} \\ \text{sc} & \textcircled{3} \\ \text{parts} & \left\langle \begin{array}{ll} \text{part} & \text{toe} \end{array} \right\rangle_{part} \\ \text{vform} & \text{fin} \\ \text{slash} & \textcircled{4} \\ \text{hebben_zijn} & \text{hebben} \end{array} \right]) \text{ :- realize-args}(\langle \textcircled{2} \mid \textcircled{1} \rangle, \textcircled{3}, \textcircled{4}).$$

Lexical Resources

- Wide-coverage of lexicalist grammars requires detailed lexical info,
- We use existing lexical resources (Celex & Parole) to obtain morphological and subcategorization info.
- Currently, the system has approx. 150K (inflected) lexical entries.

Lexical Resources

- **Celex:**

- ★ 33K lemma's for nouns, adjectives, adverbs, etc.,
- ★ 5800 lemma's for trans & intrans (particle) verbs.

- **Parole:**

- ★ 1600 verbs with subcat-frames not covered by Celex,
- ★ 800 nouns with special subcat properties.

- **“Hand”:**

- ★ 800 hand-crafted lemma's,
- ★ 4K proper names occurring in Eindhoven corpus.

Treebank

- A syntactically annotated corpus is useful for:
 - ★ Grammar Debugging,
 - ★ Evaluation,
 - ★ Collection of statistical info.
- Using current grammar directly has disadvantages:
 - ★ Grammars change,
 - ★ Annotation is difficult for strings outside coverage,
 - ★ Hard to compare with other systems,

Dependency Trees

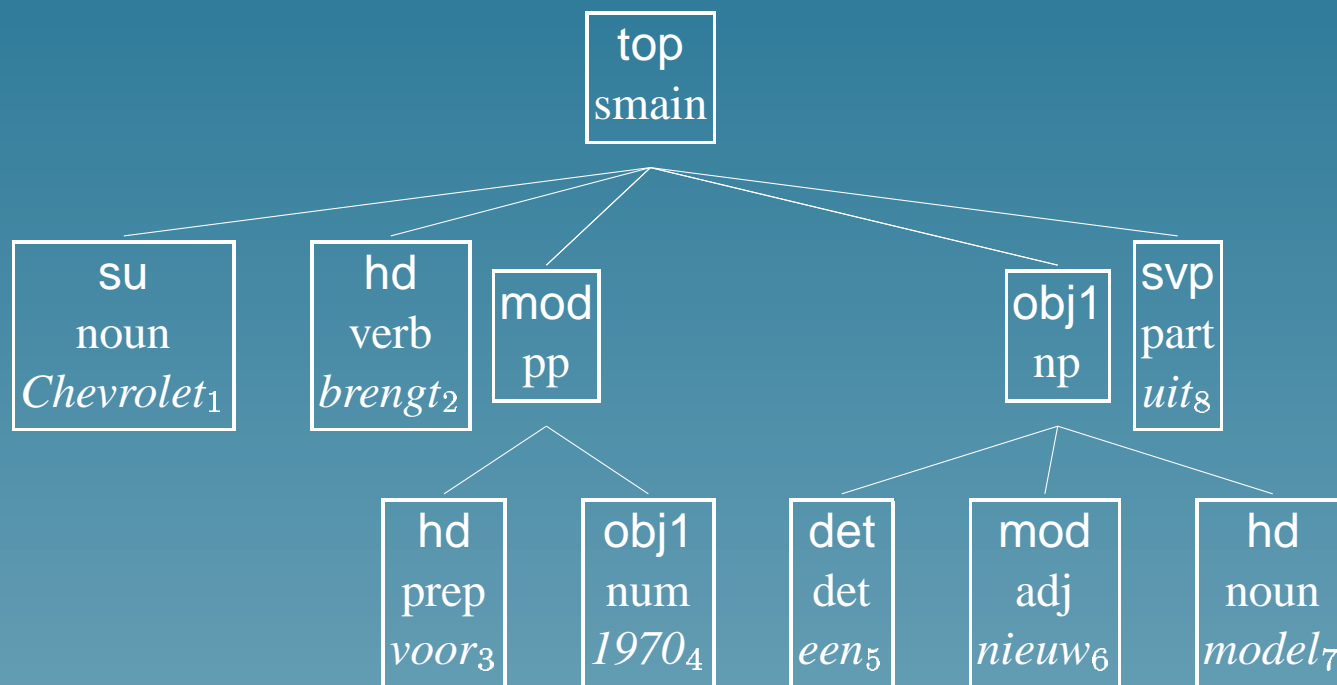
- Provide a grammar independent level of representation,
- Suitable for (relatively) free-word order languages,
- Lexical Dep Relations are useful for data-driven, statistical, parsing (Collins 98),
- We adopt annotation format for Dutch developed in CGN project.

Head-driven DT construction

- Data-structure: feature for each Dep Rel,
- A lexical head subcategorizes for a specific set of dependents, each linked to a specific Dep Rel,
- In head-comps-structures, Dep Tree can simply be shared between mother and head.

$$\text{lex}(\underset{v}{\left[\begin{array}{l} \text{deps}_{pred} \left\langle \left[\begin{array}{l} dt \\ \boxed{1} \end{array} \right], \left[\begin{array}{l} case \\ dt \end{array} \right] \right\rangle \\ \text{subj}_{np} \left[\begin{array}{l} dt \\ \boxed{3} \end{array} \right] \\ \underset{vdom}{\left[\begin{array}{ll} hwr d & vind \\ postag & verb \\ cat & inf \\ su & \boxed{3} \\ obj1 & \boxed{2} \\ predc & \boxed{1} \\ mod & \langle \rangle \end{array} \right]} \end{array} \right]}, \text{vinden}).$$

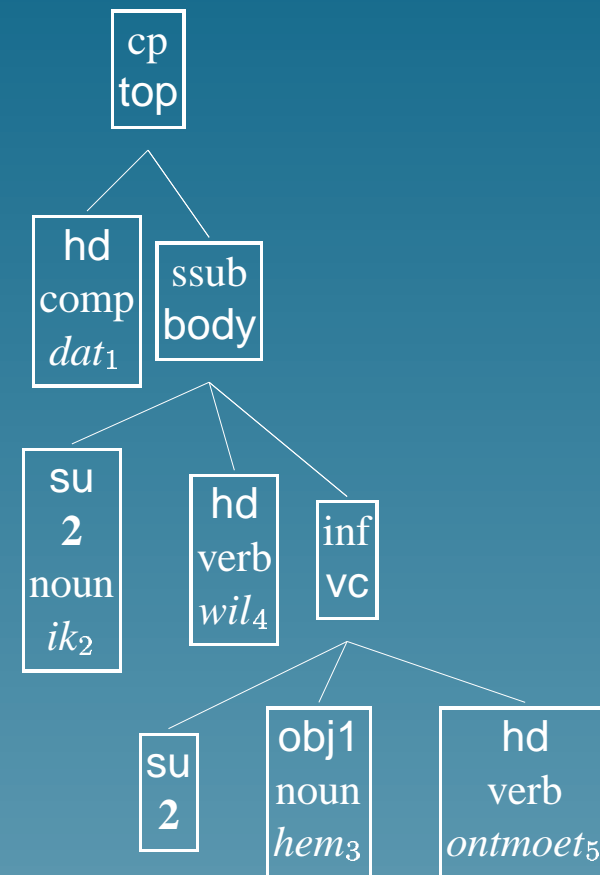
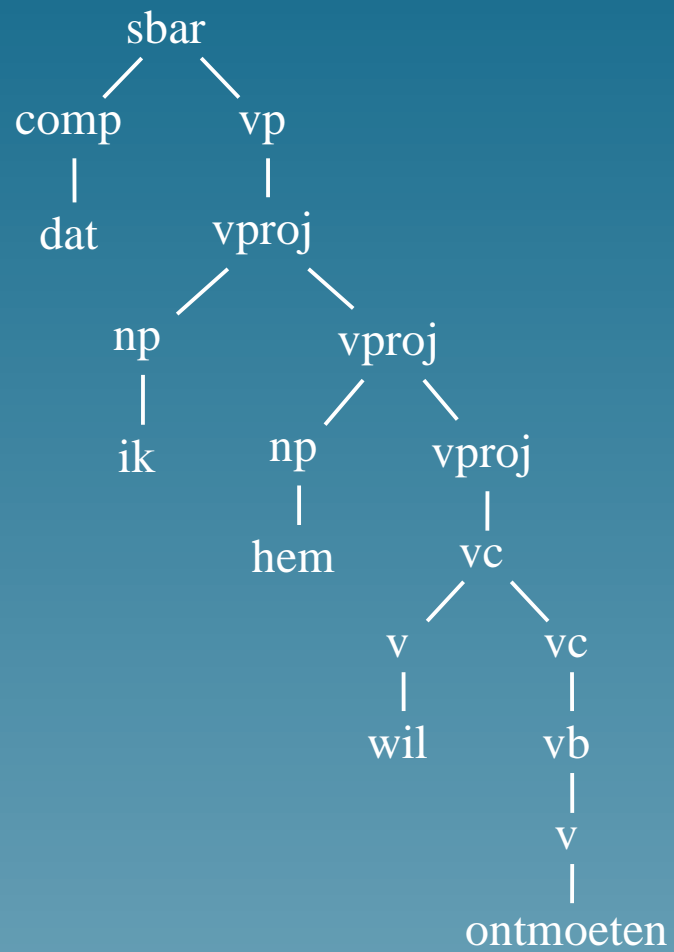
Dependency Tree



Chevrolet brengt voor 1970 een nieuw model uit :

Phrase Structure and Dep Trees

- DT-construction in the grammar:
 - ★ coordination (not a regular headed-struct),
 - ★ unbounded dependencies (not lexically headed),
 - ★ modification (no lexical treatment of adjuncts).
- Structure of Dep Tree not always isomorphic to syntactic tree.
 - ★ Example : Crossing Dependency Constructions.



Building a Treebank

- Thistle : editor for linguistic objects (Calder, 2000),
- Define a Thistle SPEC (XML DTD) for Dep Trees,
- Initial trees constructed with Alpino,
 - ★ Parse input string,
 - ★ Select (manually) best parse,
 - ★ Store corresponding Dep Tree as XML
- Use Thistle to edit and correct parse results,

Using the Treebank

- Grammar Evaluation based on Dep Rel's between lexical Heads (Carroll et al, 1999),
- Dep Tree defines as set of $\langle \text{HdWrd DepRel DepHdWrd} \rangle$, e.g.

\langle dat	body	wil	\rangle
\langle wil	su	ik	\rangle
\langle wil	vc	ontmoet	\rangle
\langle ontmoet	su	ik	\rangle
\langle ontmoet	obj1	hem	\rangle

Using the Treebank

- Parse results can be scored for precision and recall using lexically headed dependency relations,
- Useful during grammar development,
- Probabilities for lexical dependency relations can be estimated by parsing (unannotated) text,
- These can be used for disambiguation (i.e. to rank parse-results).

Conclusions

- Coverage: Combination of lexicalist HPSG-style grammar with existing lexical resources,
- Head-driven construction of Dependency Trees,
- Treebank construction,
- Grammar evaluation.

Future Work

- Expand syntactic coverage,
- Expand lexicon (use CGN lexical resources...).
- Expand treebank,
- Create parse selection tool for manual annotation,
- Build a statistical disambiguation model...