



LinGO Redwoods

— A Rich and Dynamic Treebank for HPSG —

**Stephan Oepen, Daniel P. Flickinger,
Kristina Toutanova, Christopher D. Manning**

Center for the Study of Language and Information
Stanford University

`oe@csli.stanford.edu`

Ambiguity Resolution Remains a (Major) Challenge

The Problem

- With broad-coverage grammars, even moderately complex sentences typically have multiple analyses (tens or hundreds, rarely thousands);
- unlike in grammar writing, exhaustive parsing is useless for applications;
- identifying the ‘right’ (i.e. intended) analysis is a very hard problem (AI);
- inclusion of (non-grammatical) sortal constraints is generally undesirable.

Current State of Affairs

- Heuristic scoring rules applied to (classes of) lexical items and rules;
- ‘optimality’ projection: accumulate quality marks and rank globally;
- beginning work on probabilistic models for on- or off-line parse selection.



Redwoods: Objectives and Basic Approach

Background — Motivation

- Facilitate research into stochastic disambiguation for HPSG (parsing);
- overcome limitations of existing resources (e.g. the PTB and alike);
- define and train parse ranking model over broad-coverage grammar.

A Rich and Dynamic Treebank

- Grounded in existing broad-coverage grammar and parsing system;
- hand-disambiguate: record elementary discriminating decisions;
- provide syntactic and semantic information in variable formats;
- treebank evolution and (semi-)automatic updates with the grammar.



Why (Yet) Another (Type of) Treebank?

Existing Resources (PTB, SUSANNE, NeGra, PDT, et al.)

- **(primarily) mono-stratal** topological *or* tectogrammatical;
- **(relatively) shallow** limited syntax, little or no semantics;
- **(mostly) static** (manual) ground truth annotation, no evolution.

Requirements for Disambiguation

- **syntax vs. semantics** topicalization vs. attachment ambiguity;
- **granularity** adequate match to granularity of grammar;
- **adaptability** map into various formats; semi-automated updates.



Linguistic Grammars On-Line Consortium (LinGO)

Set-Up

- Loosely organized group of institutions and interested individuals;
- (largely) shared research interests and goals, opportunistic funding;
- rooted in 'linguistic' NLP but geared towards practical applications;
- LinGO resources in wide use for research, education, applications.

History

- Originally DFKI Saarbrücken – CSLI Stanford (VerbMobil, 1994);
- Tokyo University (Tsuji Laboratory): logic compilation techniques;
- Cambridge and Sussex Universities: efficient and accurate NLP.



An Open-Source Repository: Existing Resources

Common Reference Formalism

- Strongly typed, conjunctive, closed world typed feature structure logic;
- blend of [Carpenter, 1992], [Copestake, 1992], and [Krieger, 1995].

Engineering and Processing Environments

- LKB: grammar development environment (Lisp) [Copestake, 2002];
- PET: efficient, industry quality runtime engine (C⁺⁺⁺) [Callmeier, 2000];
- [incr tsdb()]: competence and performance profiler [Oepen, 2000].

Common Grammars on Multiple Platforms

- English (CSLI): English Resource Grammar; Dan Flickinger et al.
- Japanese (DFKI, YY Technologies): Melanie Siegel & Emily M. Bender.



LinGO English Grammar: Coverage and Size

Linguistic Coverage

- 85 % of 12,000 transcribed dialogue turns from VerbMobil domains;
- average 9-word utterances, ranging from 1 – 40 words in length;
- 80 % of phenomena-based examples in HP-derived test suite.

Size of Grammar (as of July 2001)

- 8082 types (5552 leaf types) for lexicon, rules, and semantics;
- 6897 lexical entry stems (corresponds to 17917 inflected forms);
- 29 lexical (6 inflectional) and 45 phrase structure rules;
- 24,000 source lines for type definitions (excluding lexical entries);
- 30,000 lines for hand-built lexicon (more recently as relational DB).



Sample Dialogue (from VerbMobil domain) Analyzed by LinGO English Grammar

- *Are you free in the afternoon of Thursday the first, or the morning of Friday the second? Those are pretty open days for me.*
- *Well, on the first I can only meet you after four pm, so that might be tight. On Friday I have a seminar from ten to two thirty.*
- *Are you going to be around anytime on the fifth? I am free all day on the fifth.*
- *Well, on the fifth I have a seminar from ten to two and a lecture from three to five, so I could, from nine o'clock to ten in the morning, if that is not too early for you?*
- *That is perfectly fine with me, if it is okay with you.*



LinGO Redwoods: a Rich and Dynamic Treebank

- Tie treebank development to existing broad-coverage grammar;
- hand-select (or reject) intended analyses from parsed corpus;
- [Carter, 1997]: annotation by basic discriminating properties;
- record annotator decisions (and entailment) as first-class data;
- provide toolkits for dynamic mappings into various formats;
- semi-automatically update treebank as the grammar evolves;
- integrate treebank maintenance with grammar regression testing.



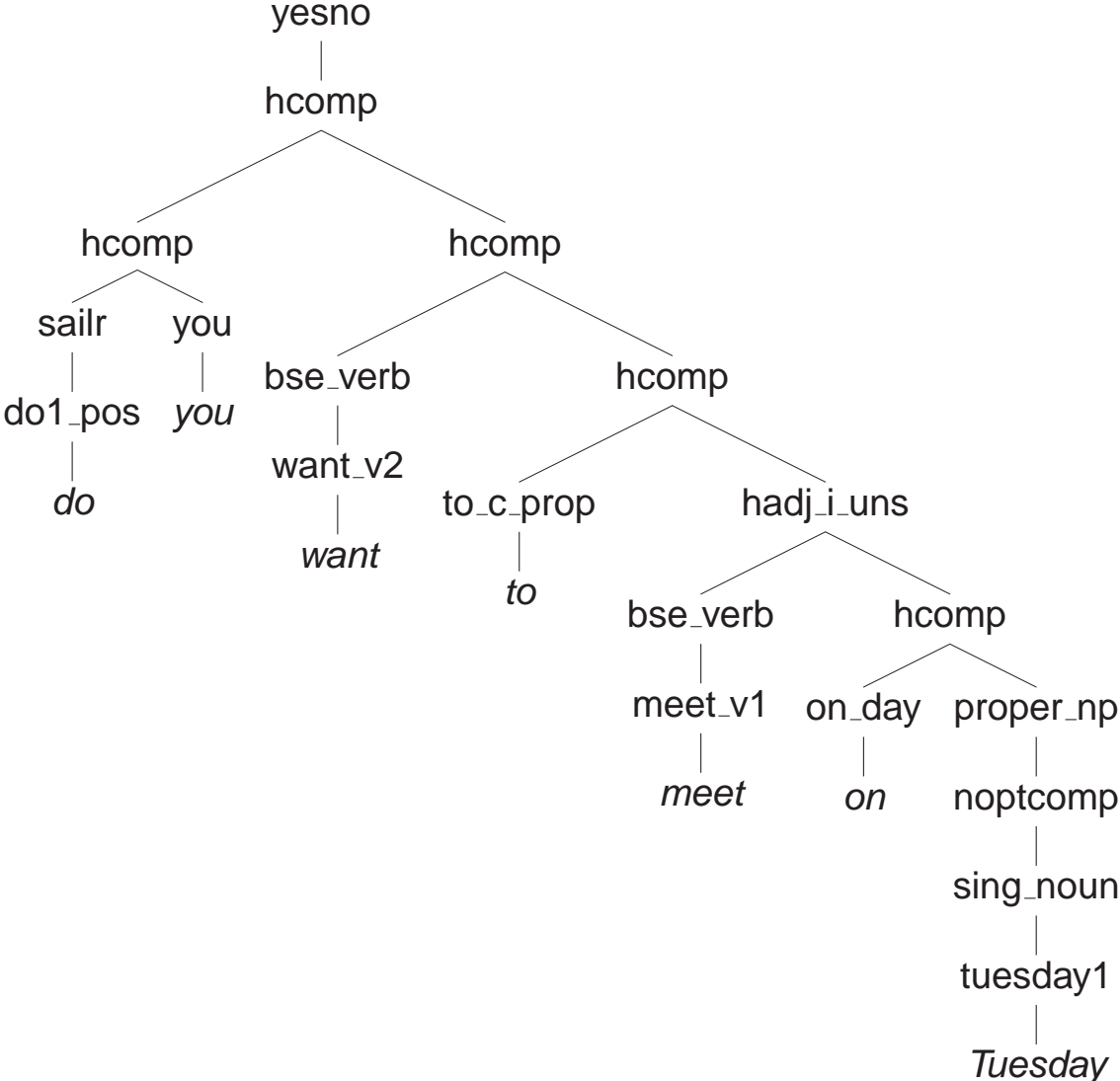
Annotation: Basic Discriminating Properties

- Extract minimal set of *basic discriminants* from set of HPSG analyses;
- typically easy to judge, need little expert knowledge about grammar;
- allow quick navigation through parse forest and incremental reduction;
- *constituents* use of particular construction over substring of input;
- *lexical items* use of particular lexical entry for input token;
- *labeling* assignment of particular abbreviatory label to a constituent;
- *semantics* appearance of particular key relation on constituent;
- Stanford undergraduate annotates some 2000 sentences per week.

- Regularly propagate discriminants into new version of parsed corpus;

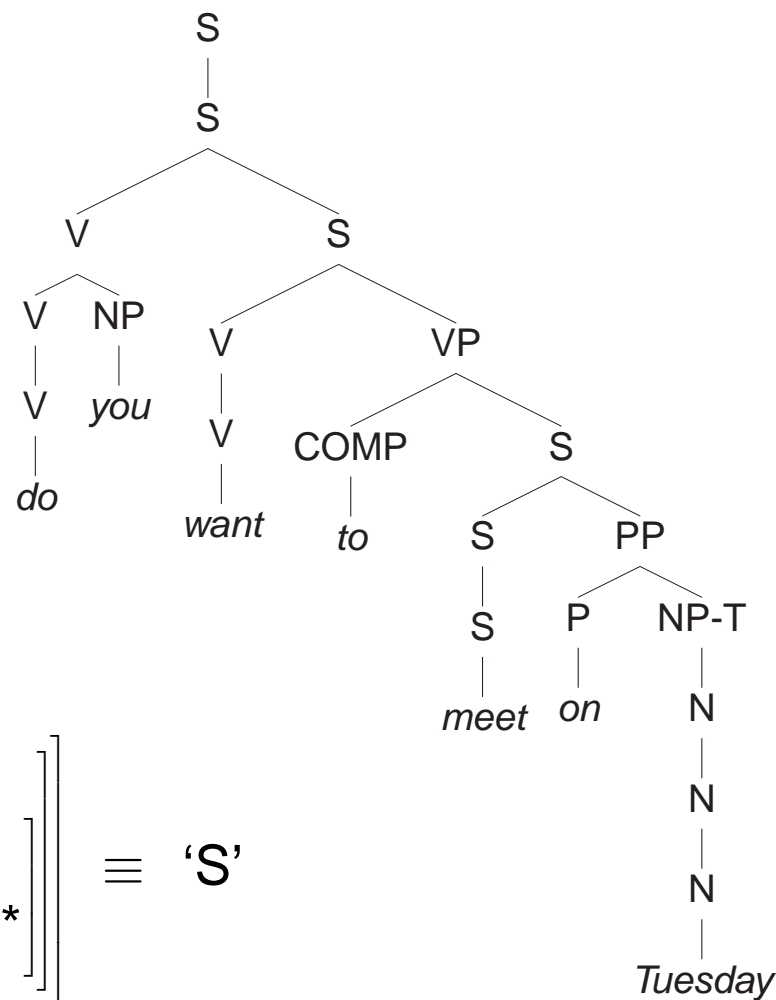


Redwoods Representations: Native Encoding



Derived Encodings: Labeled Phrase Structure Trees

- reconstruct full HPSG analysis from derivation tree;
- match underspecified feature structure 'templates' against each node;
- optionally, collapse or suppress nodes.



$$\text{label} \left[\text{SYNSEM.LOCAL.CAT} \left[\begin{array}{l} \text{HEAD } \textit{verbal} \\ \text{VAL} \left[\begin{array}{l} \text{SUBJ } \langle \rangle \\ \text{COMPS } *olist* \end{array} \right] \end{array} \right] \right] \equiv \text{'S'}$$



Derived Encodings: Elementary Dependencies

- reconstruct full HPSG analysis, compute MRS meaning representation;
 - extract basic predicate – argument structure with uninterpreted roles;
- labeled dependency graph fragments of (primarily) lexical relations.

```
_4:{  
  _4:int_rel[SOA e2:_want2_rel]  
  e2:_want2_rel[ARG1 x4:pron_rel, ARG4 _2:hypo_rel]  
  _1:def_rel[BV x4:pron_rel]  
  _2:hypo_rel[SOA e18:_meet_v_rel]  
  e18:_meet_v_rel[ARG1 x4:pron_rel]  
  e19:_on_temp_rel[ARG e18:_meet_v_rel, ARG3 x21:dofw_rel]  
  x21:dofw_rel[NAMED :tue]  
  _3:def_np_rel[BV x21:dofw_rel]  
}
```



Redwoods Third Growth: Current Development Status

	all parses			active = 0			active = 1			active > 1		
	#		×	#		×	#		×	#		×
VM₆	2706	7.7	46.7	216	9.4	63.5	2482	8.3	43.5	6	15.8	757.8
VM₁₃	2279	8.5	61.9	248	10.8	80.5	2029	8.7	59.5	2	15.5	198.0
VM₃₁	1967	6.2	27.9	216	10.1	95.9	1746	7.5	30.8	5	8.4	20.8
VM₃₂	699	7.5	53.2	15	11.8	57.7	684	8.4	53.2	0	0.0	0.0
Total	7651	7.5	47.0	695	10.2	79.5	6941	8.2	45.9	13	12.9	388.2

‘#’ total number of items (primarily sentences) in each aggregate;

‘||’ average length of items in aggregate (number of input tokens);

‘×’ average structural ambiguity (analyses assigned by the grammar).

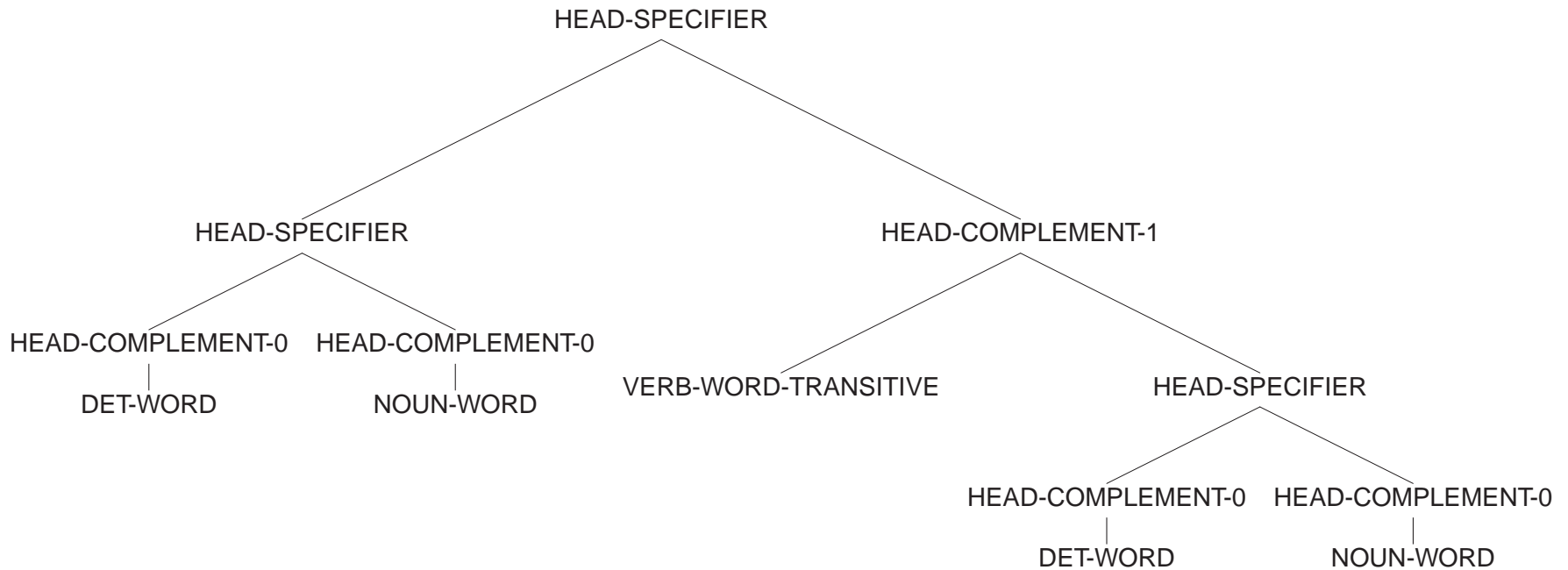


Redwoods Applications: Parse Disambiguation

- Establish experimental baseline for (simple) probabilistic models;
- restrict to Redwoods subset of fully disambiguated ambiguous items;
- HMM: maximize likelihood for sequences of preterminal types;
- PCFG: estimate rule probability from relative frequency of local tree;
- optionally condition on parent in phrase structure tree (grandparent);
- ten-fold cross validation: score against annotated gold standard;
- more recently: richer, conditional models including semantic features.



(Simplified) PCFG Estimation Example



probability	context-free rule (i.e. local tree of depth one)
$2/3 = 0.66$	HEAD-SPECIFIER \rightarrow HEAD-COMPLEMENT-0 HEAD-COMPLEMENT-0
$1/3 = 0.33$	HEAD-SPECIFIER \rightarrow HEAD-SPECIFIER HEAD-COMPLEMENT-1
$1/1 = 1.00$	HEAD-COMPLEMENT-1 \rightarrow VERB-WORD-TRANSITIVE HEAD-SPECIFIER
$2/4 = 0.50$	HEAD-COMPLEMENT-0 \rightarrow NOUN-WORD
$2/4 = 0.50$	HEAD-COMPLEMENT-0 \rightarrow DET-WORD



Redwoods Applications: Preliminary Results

Method		Tag Selection	Parse Selection
Random		90.11%	25.98%
HMM	unigram	96.51%	43.74%
	trigram	97.73%	47.37%
	perfect	100.00%	54.88%
PCFG	simple	97.20%	65.49%
	grandparent	97.52%	72.15%
	combined	98.00%	76.67%

$$combined(t) = \log(P_{grandparent}(t)) + \lambda \log(P_{trigram}(tags(t)))$$



Related Work

Non-Public Environments

- Related work at SRI Cambridge, (Xerox) PARC, and M\$ Research;
- grammars, language corpora, and treebanks not publicly available;
- results published in some cases, generally difficult to reproduce.

Academic Environments

- [Dipper, 2000] LFG for German, ‘transfer’ into NeGra format;
- [Bouma et al., 2001] HPSG for Dutch, dependency structures only;
- emerging ‘HPSG’ treebanks for Bulgarian, Polish, maybe others;
- to our best knowledge, no existing *rich* and *dynamic* approach.



Conclusions — Background Material

- ‘Deep’ grammar-based processing requires adequate stochastic models;
- basic research needed on acquisition and application of stochastic models;
- no existing treebank resources with suitable granularity and flexibility;
- LinGO Redwoods treebank based on existing open-source technology;
- tied to broad-coverage HPSG grammar: advantages and disadvantages;
- rich in available information, dynamic in data extraction and evolution.

Grammar and Treebank available from: <http://redwoods.stanford.edu/>



Outlook: Redwoods for Everyone



Towards a common, open-source basis for stochastic HPSG



Based on Research and Contributions of

Tim Baldwin, John Beavers, Ezra Callahan
Emily M. Bender, Kathryn Campbell-Kibler,
John Carroll, Ann Copestake,
Rob Malouf, Ivan A. Sag,
Stuart Shieber, Tom Wasow,
and others.