

Chapter 29

Computational linguistics and grammar engineering

Emily M. Bender

University of Washington

Guy Emerson

University of Cambridge

We discuss the relevance of HPSG for computational linguistics, and the relevance of computational linguistics for HPSG.

1 Introduction

From the inception of HPSG in the 1980s, there has been a close integration between theoretical and computational work (for an overview, see Flickinger, Polard & Wasow 2018, Chapter 2 of this volume). In this chapter, we discuss computational work in HPSG, starting with the infrastructure that supports it (both theoretical and practical) in Section 2. Next we describe several existing large-scale projects which build HPSG or HPSG-inspired grammars (see Section 3) and the deployment of such grammars in applications including both those within linguistic research and otherwise (see Section 4). Finally, we turn to linguistic insights gleaned from broad-coverage grammar development (see Section 5).

2 Infrastructure

2.1 Theoretical considerations

There are several properties of HPSG as a theory that make it well-suited to computational implementation. First, the theory is kept separate from the formalism:



Emily M. Bender & Guy Emerson. 2018. Computational linguistics and grammar engineering. In Stefan Müller, Anne Abeillé, Robert D. Borsley & Jean-Pierre Koenig (eds.), *Head-Driven Phrase Structure Grammar: The handbook*, 477–488. Berlin: Language Science Press. DOI:??

the formalism is expressive enough to encode a wide variety of possible theories. While some theoretical work does argue for or against the necessity of particular formal devices (e.g. the shuffle operator (**FIXME-Reape**)), much of it proceeds within shared assumptions about the formalism. This is in contrast to work in the context of the Minimalist Program (Chomsky 1993), where theoretical results are typically couched in terms of modifications to the formalism itself. From a computational point of view, the benefit of differentiating between theory and formalism is that it means that the formalism is relatively stable. That in turns enables the development and maintenance of software systems that target the formalism, for parsing, generation, and grammar exploration (see Section 2.3 below for some examples).¹

A second important property of HPSG that supports a strong connection between theoretical and computational work is an interest in both so-called ‘core’ and so-called ‘peripheral’ phenomena. Most implemented grammars are built with the goal of handling naturally occurring text.² This means that they will need to handle a wide variety of linguistic phenomena not always treated in theoretical syntactic work (**FIXME-Baldwin-et-al-Beauty**). A syntactic framework that excludes research on ‘peripheral’ phenomena as uninteresting provides less support for implementational work than does one, like HPSG or Construction Grammar, that values such topics (for a comparison of HPSG and Construction Grammar, see Müller 2018, Chapter 37 of this volume).

Finally, the type hierarchy characteristic of HPSG lends itself well to developing broad-coverage grammars which are maintainable over time (**FIXME-find-cite?**)³ The use of the type hierarchy to manage complexity at scale comes out of the work of Flickinger (1987) and others at HP labs in the project where HPSG was originally developed. The core idea is that any given constraint is (ideally) expressed only once on types which serve as supertypes to all entities that bear that constraint.³ Such constraints might represent broad generalizations that apply to many entities or relatively narrow, indiosyncratic properties. By isolating any given constraint on one type (as opposed to repeating it in mutiple places), we build grammars that are easier to update and adapt in light of new data that

¹There are implementations of Minimalism, notably **FIXME-Stabler** and **FIXME-Indianadiss**. However, writing an implementation requires fixing the formalism, and so these are unlikely to be useful for testing theoretical ideas as the theory moves on.

²It is possible, but less common, to do implementation work strictly against testsuites of sentences constructed specifically to focus on phenomena of interest.

³Originally this only applied to lexical entries in Flickinger’s work. Now it also applies phrase structure rules, lexical rules, and types below the level of the sign which are used in the definition of all of these.

require refinements to constraints. Having a single locus for each constraint also makes the types a very useful target for documentation (FIXME:LTDB) and grammar exploration (FIXME:typediff).

2.2 Practical considerations

The formalism of HPSG allows practical implementations, since feature structures are well-defined data structures. Furthermore, because HPSG is defined to be bi-directional, an implemented grammar can be used for both parsing and generation. In this section, we discuss how HPSG allows tractable algorithms, which enables linguists to empirically test hypotheses, and which also enables HPSG grammars to be used in a range of applications, as we will see in Sections 4.1 and 4.2, respectively.

2.2.1 Computational complexity

One way to measure how easy or difficult it is to use a syntactic theory in practical computational applications is to consider the *computational complexity* of parsing and generation algorithms (Gazdar & Pullum 1985). For example, we can consider how much computational time a parsing algorithm needs to process a particular sentence. For longer sentences, we would expect the amount of time to increase, but the more complex the algorithm is, the more quickly the amount of time increases. If we consider sentences containing n tokens, we can find the average amount of time taken, or the longest amount of time taken. We can then increase n , and see how the amount of time changes, both in the average case, and in the worst case.

At first sight, analysing computational complexity would seem to paint HPSG in a bad light, because the formalism allows us to write grammars which can be arbitrarily complex; in technical terminology, the formalism is *Turing-complete* (Johnson 1988: Section 3.4). However, as discussed in the previous section, there is a clear distinction between theory and formalism. Although the feature-structure formalism rules out the possibility of efficient algorithms that could cope with any possible feature-structure grammar, a particular theory (or a particular grammar) might well allow efficient algorithms.

The difference between theory and formalism becomes clear when comparing HPSG to other computationally-friendly frameworks, such as Combinatory Categorical Grammar (CCG),⁴ or Tree Adjoining Grammar (TAG; Joshi 1987; Schabes

⁴For an introduction, see Steedman & Baldridge (2011). For a comparison with HPSG, see Kubota (2018), Chapter 34 of this volume.

et al. 1988)). The formalisms of CCG and TAG inherently limit computational complexity: for both of them, as the sentence length n increases, worst-case parsing time is proportional to n^6 (Kasami et al. 1989). This is a deliberate feature of these formalisms, which aim to be just expressive enough to capture human language, and not any more expressive. Building this kind of constraint into the formalism itself highlights a different school of thought from HPSG. Indeed, Müller (2015) explicitly argues in favor of developing linguistic analyses first, and improving processing efficiency second. As discussed above in Section 2.1, separating the formalism from the theory means that the formalism is stable, even as the theory develops.

It would be beyond the scope of this chapter to give a full review of parsing algorithms, but it is instructive to give an example. For grammars that have a context-free backbone (we can express every derivation as a phrase-structure tree plus constraints between mother and daughter nodes), it is possible to adapt the standard chart-parsing algorithm for context-free grammars. The basic idea is to parse “bottom-up” through the tree, starting by finding analyses for each token in the input, and then finding analyses for increasingly longer sequences of tokens, until we reach the entire sentence. The resulting algorithm is more computationally complex than for a context-free grammar, because we are dealing with feature structures, rather than nonterminal symbols. While a context-free grammar allows a finite number of nonterminals, a feature-structure grammar may allow a infinite number of possible feature structures. For an HPSG grammar without recursive unary rules, this algorithm has a worst-case complexity of exponential time. This is less complex than for an arbitrary grammar (which means that this class of grammars is *not* Turing-complete), but more complex than for CCG or TAG. However, when parsing real corpora, it turns out that the average-case complexity is much better than we might expect (Carroll 1994). On the one hand, grammatical constructions do not generally combine in the worst-case way, and on the other hand, when a grammar writer is confronted with multiple possible analyses for a particular construction, they may opt for the analysis that is more efficient for a particular parsing algorithm. Oepen & Flickinger (1998) describe the use of test-suites for profiling the time efficiency of a grammar and parsing algorithm, when running on a representative sample of text.

2.2.2 Parse ranking

For an ambiguous sentence, a grammar gives multiple valid parses. This is widely known, with attachment ambiguities and coordination ambiguities being partic-

ularly well-known examples. However, people are naturally very good at resolving ambiguity, which means most ambiguity is not apparent, even to linguists. It is only with the development of large-scale grammars that the sheer scale of ambiguity has become clear. For example, (1) might seem unambiguous, but there is a second reading, where *my favorite* is the topicalized object of *speak*, which would mean that town criers generally speak the speaker's favorite thing (perhaps a language) clearly. There is also a third, even more implausible reading, where *my favorite town* is the topicalized object. Such implausible readings don't easily come to mind, and in fact, the 2018 version of the English Resource Grammar (Flickinger 2000; 2011) gives a total of 21 readings for this sentence. With increasingly long sentences, such ambiguities stack up very quickly. For (2), the first line of a newspaper article,⁵ the ERG gives 35,094 readings.

- (1) My favorite town criers speak clearly.
- (2) A small piece of bone found in a cave in Siberia has been identified as the remnant of a child whose mother was a Neanderthal and father was a Denisovan, a mysterious human ancestor that lived in the region.

While exploring ambiguity can be interesting for a linguist, typical practical applications require just one parse per input sentence and specifically the parse the best reflects the intended meaning (or only the top few parses, in case the one put forward as 'best' might be wrong). Thus, what is required is a *ranking* of the parses, so that the application can only use the most highly-ranked parse, or the top *N* parses.

Parse ranking is not usually determined by the grammar itself, because of the difficulty of manually writing disambiguation rules. Typically, a statistical system is used (Toutanova et al. 2002; 2005). First, a corpus is *treebanked*: for each sentence in the corpus, an annotator (often the grammar writer) chooses the best parse, out of all parses produced by the grammar. The set of all parses for a sentence is often referred to as the *parse forest*, and the selected best parse is often referred to as the *gold standard*. Given the gold parses for the whole corpus, a statistical system is trained to predict the gold parse from a parse forest, based on many features⁶ of the parse. From the example in (1), we can see how a number of a different features all influence the preferred interpretation: the likelihood of a construction (such as topicalization), the likelihood of a valence frame (such

⁵<https://www.theguardian.com/science/2018/aug/22/offspring-of-neanderthal-and-denisovan-identified-for-first-time>

⁶In the machine-learning sense of "feature", not the feature-structure sense.

as transitive *speak*), the likelihood of a collocation (such as *town crier*), the likelihood of a semantic relation (such as speaking a town), and so on.

Because of the large number of possible parses, it can be helpful to *prune* the search space: rather than ranking the full set of parses, we can restrict attention to a smaller set of parses, which hopefully includes the correct parse. By carefully choosing how to restrict our attention, we can drastically reduce processing time without hurting parsing accuracy. One method, called *supertagging*,⁷ exploits the fact that HPSG is a lexicalized theory: choosing the correct lexical entry brings in rich information that can be exploited to rule out many possible parses. Thus if the correct lexical entry can be chosen prior to parsing (e.g. on the basis of the prior and following words), the range of possible analyses the parser must consider is drastically reduced. Although there is a chance that the supertagger will predict the wrong lexical entry, using a supertagger can often improve parsing accuracy, by ruling out parses that the parse-ranking model might incorrectly rank too high. Supertagging was first applied to HPSG by Matsuzaki et al. (2007), building on previous work for TAG (Bangalore & Joshi 1999) and CCG (Clark & Curran 2004). To allow multi-word expressions (such as *by and large*), where the grammar assigns a single lexical entry to multiple tokens, Dridan (2013) has proposed an extension of supertagging, called *ubertagging*, which jointly predicts both a segmentation of the input and supertags for those segments. Dridan manages to increase parsing speed by a factor of four, while also improving parsing accuracy.

Finally, in order to train these statistical systems, we need to first annotate a treebank. When there are many parses for a sentence, it can be time-consuming to select the best parse. To efficiently use an annotator’s time, it can be helpful to use *discriminants*, properties which hold for some parses but not for others (Carter 1997). For example, discriminants might include whether to analyse an ambiguous token as a noun or a verb, or where to attach a prepositional phrase. This approach to treebanking also means that annotations can be re-used when the grammar is updated (OFTM2004a-u; Flickinger et al. 2017).

2.2.3 Semantic dependencies

In practical applications of HPSG grammars, the full derivation trees and the full feature structures are often unwieldy, containing far more information than

⁷The name refers to *part-of-speech tagging*, which predicts a part-of-speech for each input token, from a relatively small set of part-of-speech tags. Supertagging is “super”, in that it predicts detailed lexical entries, rather than simple tags.

necessary for the task at hand. It is therefore often desirable to extract a concise semantic representation.

In computational linguistics, a popular approach to semantics is to represent the meaning of a sentence as a *dependency graph*, as this enables the use of graph-based algorithms.⁸ Several types of dependency graph have been proposed based on Minimal Recursion Semantics (MRS; Copestake et al. 2005), with varying levels of simplification. The most expressive is Dependency Minimal Recursion Semantics (DMRS; Copestake 2009), which is fully interconvertible with MRS.⁹ In contrast, Elementary Dependency Structures (EDS; Oepen & Lønning 2006) lose some scope information, which, for many applications, is less important than predicate-argument structure. Finally, DELPH-IN MRS Dependencies (DM; Ivanova et al. 2012) express predicate-argument structure purely in terms of the surface tokens, without introducing any abstract predicates. A comparison of MRS, DMRS, and DM is given in Figures 1 and 2. The existence of such dependency graph formalisms, as well as software packages to manipulate such graphs (`copestake2016pydmrs`), has made it easier to use HPSG grammars in a number of practical tasks, as we will discuss in Section 4.2.

2.3 A brief history of HPSG grammar engineering

History: PAGE, VerbMobil, ??

Current platforms:

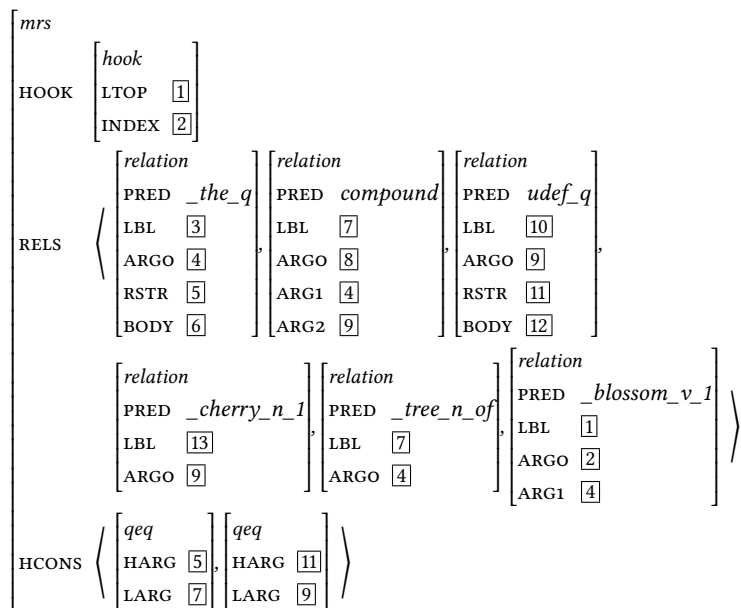
- LKB/ACE/PET/Agree
- Trale
- Other

3 Development of HPSG resources

- CoreGram
- DELPH-IN consortium

⁸In this section, we are concerned with *semantic* dependencies. For *syntactic* dependencies, see Hudson (2018), Chapter 36 of this volume.

⁹More precisely, there is a one-to-one correspondence between DMRS and MRS structures, if every predicate is a unique *intrinsic variable*. As observed by Oepen & Lønning (2006), this allows a variable-free semantic representation, by replacing each reference to a variable with a reference to the corresponding predicate.



(a) MRS, as a feature structure.

INDEX : e_1

$l_1 : \text{_the_q}(x_1, h_1, h_2), h_1 \text{ QEQ } l_4$

$l_2 : \text{udef_q}(x_2, h_3, h_4), h_3 \text{ QEQ } l_3$

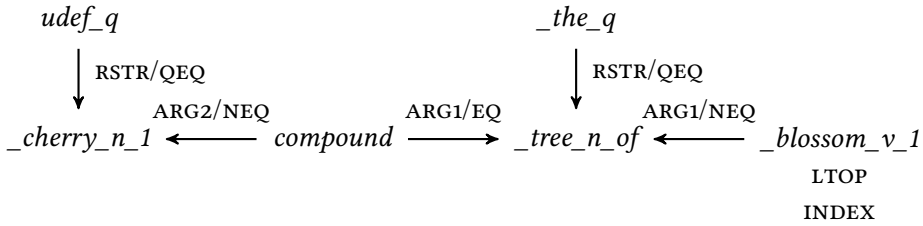
$l_3 : \text{_cherry_n_1}(x_2)$

$l_4 : \text{_tree_n_of}(x_1), \text{compound}(e_2, x_1, x_2)$

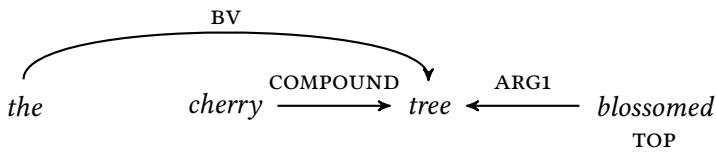
LTOP $l_5 : \text{_blossom_v_1}(e_1, x_1)$

(b) MRS, abstractly.

Figure 1: Minimal Recursion Semantics (MRS) representation for the sentence “the cherry tree blossomed”, as produced by the English Resource Grammar (ERG). For simplicity, we have omitted some details, including features such as number and tense, individual constraints (ICONS), and the use of difference lists.



(a) Dependency Minimal Recursion Semantics (DMRS) is as expressive as MRS. Predicates are represented as nodes, while semantic roles and scopal constraints are represented as links.



(b) DELPH-IN MRS Dependencies (DM) are a simplified version of MRS. Nodes are the tokens in the sentence (rather than abstract predicates), and some scopal information is dropped.

Figure 2: MRS-based dependency graphs for the sentence “the cherry tree blossomed”, based on the MRS given in Figure 1.

- ERG
- Other large-ish grammars
- Grammar Matrix
- Systems inspired by HPSG:
 - Alpino
 - Enju

4 Deployment of HPSG resources

4.1 Language documentation and linguistic hypothesis testing

Deployment for linguistic goals.

4.1.1 CoreGram

4.1.2 Grammar Matrix

4.1.3 AGGREGATION

4.1.4 Derived resources: Redwoods-style treebanks

4.2 Downstream applications

In this section, we discuss the use of HPSG grammars for practical tasks. There is a large number of applications, and we focus on several important applications here. Within the DELPH-IN community, a regularly updated list of applications is available online.¹⁰

4.2.1 Education

Precise syntactic analyses can be useful in language teaching, in order to automatically identify errors and give feedback to the student. In order to model common mistakes, a grammar can be extended with so-called “mal-rules”. A mal-rule is like a normal rule, in that it licenses a construction, and can be treated the same during parsing – however, given a parse, the presence of mal-rules indicates that the student needs to be given feedback (**Ben:Fli:Oep:04**; **flickinger2013error**; **morgadodacosta2016error**). A large scale system implementing this kind of computer-aided teaching has been developed by the Education Program for Gifted Youth at Stanford University, using the ERG (**suppes2014teach**). This system has reached tens of thousands of elementary and middle school children, and has been found to improve the school results of underachieving children.

Another way to use a precision grammar is to automatically produce teaching materials. Given a semantic representation, a grammar can generate one or more sentences. **Flickinger:17** uses the ERG to produce practice exercises for a student learning first-order logic. For each exercise, the student is presented with a natural language sentence, and is supposed to write down the corresponding first-order logical form. By using a grammar, the system can produce syntactically varied questions, and automatically evaluate the student’s answer.

4.2.2 NLP tasks

HPSG grammars have been used in a range of tasks in Natural Language Processing (NLP).

¹⁰<http://moin.delph-in.net/DelphinApplications>

Information extraction. **schaefer2011acl reiplinger2012glossary**

Summarization. **fang2016summarise**

Machine translation. Parse, transfer, generate. (**OVL2007a-u; bond2011deep**).
goodman2018translate Statistical without transfer, **horvat2017translate**

Distributional semantics. **emerson2018functional**

4.2.3 Data for machine learning

In section 4.2.2, we described how HPSG grammars can be used to tackle a number of NLP tasks. Another use of HPSG grammars in NLP is to generate data, on which a statistical system can be trained.

Training deep learning systems – semantic parsing – skip over HPSG, go straight to semantic representations. (**oepen2014semeval; oepen2015semeval**).
buys2017parse chen2018parse

Evaluation of deep learning systems – ShapeWorld – use a grammar to produce annotations. **kuhnle2018shapeworld**

4.3 Other?

Alpino?

Enju?

5 Linguistic Insights

- Ambiguity
- Long-tail phenomena (raising and control?)
- Scaling up (thematic roles)
- CLIMB methodology

6 Summary

Abbreviations

Acknowledgements

References

- Bangalore, Srinivas & Aravind K. Joshi. 1999. Supertagging: an approach to almost parsing. *Computational Linguistics* 25(2). 237–265. <http://aclweb.org/anthology/J99-2004>.
- Carroll, John. 1994. Relating complexity to practical performance in parsing with wide-coverage unification grammars. In James Pustejovsky (ed.), *32th Annual Meeting of the Association for Computational Linguistics. Proceedings of the conference*, 287–294. Las Cruces: Association for Computational Linguistics.
- Carter, David. 1997. The TreeBanker. A tool for supervised training of parsed corpora. In *Proceedings of the 1997 ACL workshop on computational environments for grammar development and linguistic engineering (ENVGRAM)*, 9–15. Madrid, Spain. <http://aclweb.org/anthology/W97-1502>.
- Chomsky, Noam. 1993. A Minimalist Program for linguistic theory. In Kenneth Hale & Samuel Jay Keyser (eds.), *The view from building 20: Essays in linguistics in honor of Sylvain Bromberger* (Current Studies in Linguistics 24), 1–52. Cambridge, MA/London: MIT Press.
- Clark, Stephen & James R. Curran. 2004. The importance of supertagging for wide-coverage CCG parsing. In *Proceedings of the 20th international conference on computational linguistics (COLING)*. <http://aclweb.org/anthology/C04-1041>.
- Copestake, Ann. 2009. Slacker semantics: why superficiality, dependency and avoidance of commitment can be the right way to go. In *Proceedings of the 12th conference of the European chapter of the Association for Computational Linguistics (EACL)*, 1–9. <http://aclweb.org/anthology/E09-1001>.
- Copestake, Ann, Daniel P. Flickinger, Carl J. Pollard & Ivan A. Sag. 2005. Minimal Recursion Semantics: An introduction. *Research on Language and Computation* 3(2–3). 281–332. DOI:10.1007/s11168-006-6327-9
- Dridan, Rebecca. 2013. Ubertagging: joint segmentation and supertagging for English. In *Proceedings of the 2013 conference on empirical methods in natural language processing (EMNLP)*, 1201–1212. <http://aclweb.org/anthology/D13-1120>.

- Flickinger, Dan, Stephan Oepen & Emily M. Bender. 2017. Sustainable Development and Refinement of Complex Linguistic Annotations at Scale. In Nancy Ide & James Pustejovsky (eds.), *Handbook of Linguistic Annotation*, 353–377. Dordrecht: Springer Netherlands. http://dx.doi.org/10.1007/978-94-024-0881-2_14. DOI:10.1007/978-94-024-0881-2_14
- Flickinger, Daniel P. 1987. *Lexical rules in the hierarchical lexicon*. Stanford University dissertation.
- Flickinger, Daniel P. 2000. On building a more efficient grammar by exploiting types. *Natural Language Engineering* 6(1). 15–28.
- Flickinger, Daniel P. 2011. Accuracy vs. robustness in grammar engineering. In Emily M. Bender & Jennifer E. Arnold (eds.), *Language from a cognitive perspective: Grammar, usage, and processing*, 31–50. Stanford, CA: CSLI Publications.
- Gazdar, Gerald & Geoffrey K. Pullum. 1985. Computationally relevant properties of natural languages and their grammars. *New Generation Computing* (3). 237–306. Reprinted in Savitch, Bach, Marsh, and Safran-Naveh (eds.), *The Formal Complexity of Natural Language*.
- Ivanova, Angelina, Stephan Oepen, Lilja Øvrelid & Dan Flickinger. 2012. Who did what to whom?: A contrastive study of syntacto-semantic dependencies. In *Proceedings of the 6th linguistic annotation workshop*, 2–11. <http://aclweb.org/anthology/W12-3602>.
- Johnson, Mark. 1988. *Attribute-value logic and the theory of grammar* (CSLI Lecture Notes 16). Stanford, CA: CSLI Publications.
- Joshi, Aravind K. 1987. Introduction to Tree Adjoining Grammar. In Alexis Manaster Ramer (ed.), *The mathematics of language*, 87–114. Amsterdam: John Benjamins Publishing Co.
- Kasami, Tadao, Hiroyuki Seki & Mamoru Fujii. 1989. Generalized context-free grammars and multiple context-free grammars. *Systems and Computers in Japan* 20(7). 43–52.
- Matsuzaki, Takuya, Yusuke Miyao & Jun’ichi Tsujii. 2007. Efficient HPSG parsing with supertagging and CFG-filtering. In *Proceedings of the 20th international joint conference on artificial intelligence*, 1671–1676. <https://www.aaai.org/Papers/IJCAI/2007/IJCAI07-270.pdf>.
- Müller, Stefan. 2015. The CoreGram project: Theoretical linguistics, theory development and verification. *Journal of Language Modelling* 3(1). 21–86. DOI:10.15398/jlm.v3i1.91
- Oepen, Stephan & Daniel P. Flickinger. 1998. Towards systematic grammar profiling: Test suite technology ten years after. *Journal of Computer Speech and Language* 12(4). 411–436. <http://www.delph-in.net/itsdb/publications/profiling.ps.gz>, accessed 2018-2-25. (Special Issue on Evaluation).

- Oepen, Stephan & Jan Tore Lønning. 2006. Discriminant-based MRS banking. In *Proceedings of the 5th international conference on language resources and evaluation (LREC)*, 1250–1255. http://www.lrec-conf.org/proceedings/lrec2006/pdf/364_pdf.pdf.
- Schabes, Yves, Anne Abeillé & Aravind K. Joshi. 1988. *Parsing strategies with ‘lexicalized’ grammars: Application to Tree Adjoining Grammars*. Technical Report MS-CIS-88-65. University of Pennsylvania Department of Computer & Information Science.
- Steedman, Mark & Jason Baldridge. 2011. Combinatory Categorical Grammar. In Robert D. Borsley & Kersti Börjars (eds.), *Non-transformational syntax: Formal and explicit models of grammar: A guide to current models*, 181–224. Oxford, UK/Cambridge, MA: Blackwell Publishers Ltd.
- Toutanova, Kristina, Christopher D. Manning, Dan Flickinger & Stephan Oepen. 2005. Stochastic HPSG parse disambiguation using the Redwoods corpus. *Research on Language & Computation* 3(1). 83–105.
- Toutanova, Kristina, Christopher D. Manning, Stuart M. Shieber, Dan Flickinger & Stephan Oepen. 2002. Parse disambiguation for a rich HPSG grammar. In *Proceedings of the first workshop on treebanks and linguistic theories (TLT2002)*. Sozopol, Bulgaria. <http://bultreebank.org/wp-content/uploads/2017/05/paper17.pdf>.