

NOTICE WARNING CONCERNING COPYRIGHT RESTRICTIONS:
The copyright law of the United States (title 17, U.S. Code) governs the making of photocopies or other reproductions of copyrighted material. Any copying of this document without permission of its author may be prohibited by law.

A Parser for HPSG

Alex Pranz

Jvdy 1990

Report No. CMU-LCL-90-3

**Laboratory for
Computational
Linguistics**

139 Baker Hall
Department of Philosophy
Carnegie Mellon University
Pittsburgh, PA 15213

A Parser for HPSG

Alex Franz¹

¹I would like to thank Carl Pollard and Bob Carpenter for helpful comments and discussions. This work was supported by the NSF under grant number IRI-8806913.

Contents

1 Overview	3
2 Running the Parser	3
2.1 Preparing a Grammar	3
2.2 Loading the Parser	3
2.3 Reading an Input Grammar	4
2.4 Parsing Sentences	5
2.5 Displaying Output	7
2.6 Using &TEX to Display Output	8
2.7 Changing Search Strategy and other Parameters	10
3 About HPSG	12
3.1 Organization of HPSG Grammars	12
3.2 The Signature	13
3.3 Principles of Grammar and Sort Definitions	14
3.4 The Lexicon	17
3.5 Lists	17
3.6 Relations	17
4 An Example Grammar	18
4.1 Scope of the Example Grammar	18
4.2 The Sort Signature	19
4.3 Principles of Grammar	22
4.4 Lexical Entries	27
4.5 Sort Definitions	29
5 Input Specifications	30
5.1 Signature Specification	31
5.2 Sort Definitions	31
5.3 Relation Definitions	32
6 Grammar Debugging	33
6.1 Error Messages	33
6.2 Structure of the Trace File	35
6.3 Semantic Grammar Mistakes	36
6.4 If You're Desperate	37
6.4.1 Verbose Printing of Structures	37
6.4.2 Examining the Converted Grammar	38

7 Grammar Conversion	38
7.1 Overview	38
7.2 Parsing the Signature	39
7.3 Reading Sort and Relation Definitions	40
7.4 Handling the Sort Hierarchy	40
7.5 Graph Representation	41
7.6 Constructing Initial Record Types	41
7.7 Constructing Graphs from Definitions	42
7.8 Expanding Graphs	43
8 Parsing	44
8.1 Overview	44
8.2 Initializing Search	44
8.3 Selecting Moves	45
8.4 Generating new States through Unification	46
9 Further Work	47
9.1 Efficiency Improvements	47
9.2 Extensions to the Parser	48
Appendices	50
A Sample Grammar	50
B System Trace	61
C Parser Output	65
D Printing Feature Structures	72
D.1 Overview	72
D.2 Tagging Graphs	72
D.3 Conversion to Lists	72
D.4 Printing ASCII Graphs	73
D.5 Producing MfcX Output	73
D.6 Parameters for MfcX output	74
References	78

1 Overview

This document describes a parser for Head-driven Phrase Structure Grammar (HPSG). The system was implemented in Lucid Common Lisp (Version 2.1.1) in the Laboratory for Computational Linguistics on an HP 9000/370 workstation.

Following [Carpenter and Pollard, 1989], I view the HPSG formalism as a set of constraints over a type scheme for sorted feature structures. The function of the HPSG interpreter is to solve these constraint equations.

The interpreter has three conceptual components: The grammar conversion routines transform an HPSG grammar into the appropriate internal format. The **grammar processor** performs parsing, and output routines print sorted feature structures either in ASCII format, or directly generate L^AT_EX instructions.

The next section provides a brief guide to running the parser. Subsequent sections describe HPSG from a grammar-writer's point of view, and provide formal input specifications that define the version of the HPSG formalism adopted here.

Following this, there are sections that describe grammar conversion and parsing in detail.

A sample grammar, a grammar conversion and parsing trace, some output from the program, and further details on printing feature structures are contained in the appendices.

2 Running the Parser

Following is a brief guide to running the parser that should provide enough information to allow the reader to start experimenting with the system.

2.1 Preparing a Grammar

First of all, an input grammar must be provided. Appendix A contains a sample grammar that could be modified for this purpose. Sections 4 and 5 describe HPSG grammars in more detail, and should be referred to during more serious grammar-writing efforts.

2.2 Loading the Parser

The parser has been compiled for Lucid CL versions 2.1.1 and 3.0, but the Lisp source should load in most Common Lisp implementations. When compiling the program on other machines care must be taken to first load the Lisp source files in order to properly compile macros.

The source code has been divided into the following files:

- `startup.lisp` (global declarations)
- `aux.lisp` (auxiliary functions)
- `graph-rep.lisp` (functions for representation and manipulation of graphs)
- `parse.lisp` (parser functions)
- `convert-input-gramxnar.lisp` (functions to read and convert input grammars)
- `latex-output.lisp` (functions to generate &TeX instructions for printing graphs)
- `debug.lisp` (auxiliary grammar and parser debugging functions)

Typing `(load "startup")` at the Lisp prompt causes all parser files to be loaded:

```
> (load "startup")
;;; Loading binary file "startup.hbin"
;;; Loading binary file "aux.hbin"
;;; Loading binary file "graph-rep.hbin"
;;; Loading binary file "parse.hbin"
;;; Loading binary file "convert-input-grammar.hbin"
;;; Loading binary file "latex-output.hbin"
;;; Loading binary file "debug.hbin"
#P17users/alex/Courses/Nlp3/Code/startup.hbin"
>
```

If the parser files have not been compiled for the machine in use, *first* load the source files, then type

`(compile-parser)`

to cause all the files to be compiled. Finally, load the binaries by typing `(load "startup")` again.

2.3 Reading an Input Grammar

Type `(run)` to start the parser. The greeting will be printed, and you will be asked whether you wish to keep a trace in the log file; typing `<retura>` causes this to default to "no". The trace, which is written to a file called `parser.log`, contains a lot of information (see Section 6.2) and keeping a trace will slow parsing down considerably. However, it is essential for grammar-debugging, and examining selected parts of the trace will provide insights on the working of the parser.

> (run)

pr

HPSG Parser

??? Do you want to keep a trace in the log file?: ['No'] yes

** Trace will be written to file parser.log

Next, you will be prompted to type the name of the file containing the input grammar. Type <return> to read the default file, grammarl.text:

??? Enter file name of input grammar: ["grammarl.text¹¹"]

*** Initializing Converter Variables

*** Converting grammar in file "grammarl.text¹¹"

*** Parsing signature specification

*** There are 131 sorts

** Reading sort and relation definitions

** Performing inheritance in sort hierarchy

*** Constructing minimal satisfiers

** Constructing datatypes

*** Constructing satisfiers of definite constraints

*** Constructing satisfiers of conjuncts in definitions

** Constructing satisfiers of disjunctive constraints

*** Expanding minimal satisfiers

** Grammar conversion was successful.

** Initializing parser

The parser now reads the input grammar and converts it into a number of internal data structures. Messages indicating its progress are printed throughout, as shown above.

2.4 Parsing Sentences

After a little while you will reach the parsing component of the system, indicated by the following message and prompt:

<parser> Enter type of expression and words terminated by ';>:

-->

Input to the parser consists of two parts. The first indicates the type of expression you wish to parse; possible expressions include sentence, phrase, sign. This is followed by the words that constitute the expression. Commands to the parser must be ended by a semicolon.

sentence may be abbreviated as *s*. In addition, if the keyword `:all` is included the parser will return all possible parses; otherwise, it will deliver only the first parse found. Below are some examples of possible input to the parser if the grammar in **grammar1.text** has been loaded:

```
—> sentence sandy laughs;
—> :all s sandy gives cindy tea;
—> :all phrase likes cindy;
```

Once an expression has been entered, the parser starts working on the input. It attempts to resolve sort and relation definitions in the grammar by searching through the space of non-disjunctive sorted feature structures created by disjunctions in the definitions. At every state, a message of the following format is printed out:

```
*** 0 states in queue: NIL
Moving from graph #1
Resolving disjuncts of — sentence —
Mode: DEPTH-TIL-SPLIT Generated 3 new states:
(#2 sentence #3 sentence #4 sentence)
```

The first line indicates how many unresolved states are currently waiting in the queue. Every time a new graph is created during search it is labelled with a number. The above message was printed at the initial state; thus, there are no states in the queue, and we are moving from the start state, corresponding to "graph #1". The third line tells us that the parser is attempting to resolve the definition of the sort *sentence*, a special sort created to act as a template for parsing sentences.

Finally, the last line informs us that we are in search mode DEPTH-TIL-SPLIT and that we have generated three new states. It also tells us the numbers of the new graphs (#2, #3, and #4). The term "Depth-til-split" means breadth-first search that proceeds down paths until they split. That is, at any point where only one new state is generated we continue search from the new state, but whenever two or more new states are generated the new states are adjoined to the back of the queue, and search proceeds in the usual breadth-first manner.

If all moves from some non-final state fail, a message like the following is printed:

```
⋮
** 2 states in queue: (#3 #4)
Moving from graph #2
Resolving disjuncts of — word —
All moves failed (backtracking)...
*• 1 states in queue:(#4)
Moving from graph #3
⋮
```

Graph #2, while not being final, turned out to be a dead end since a "word" constraint could not be satisfied. Thus, the graph was deleted, and search continued from the next state in the queue (#3).

If you are running in all-paths-parsing mode (having included the `:all` keyword), the parser prints this message every time it has found a parse:

```
** Found one parse — looking for more...
```

The parser may be interrupted at any time; the feature structures corresponding to the parses found so far will be accessible in a list as the value of the global variable `*all-parses*`.

Eventually, one of the following messages will be printed out:

```
*** Exiting; first parse found is in *parse*
```

```
*** Exiting; there was no successful parse
```

```
*** There are no more parses; all parses found are in *all-parses*
```

This will be followed by the parser's listener prompt

```
<parser> Enter type of expression and words terminated by ';':  
-->
```

The parser has finished, and you are ready to inspect the result.

2.5 Displaying Output

The parser listener is a simple loop that accepts expressions to be parsed (as shown above), the command `quit`, or a Lisp form to be evaluated. The latter is achieved by typing the following:

```
--> eval <Lispform>;
```

The function `(print-graph <graph>)` can be used to print a feature structure in ASCII format, as shown on the next page.

```

-> eval (print-graph "parse");
- Graph number "020" -
PHOK <<&st-pho>-word>
ir-: FIRST [1] <MBdy>
BEST W <clUt-phoB-word>
FIRST [6] <laugh>
REST <clst-phoB-word>
SYKSEM [1] <synsem>
LOC <oc> [1]
CAT <c> *t>
MARKING <<nm>rked>
BEAD [5] <*>
VPOBif <fin>
AUX <minw>
IKV <nun>
PRD <mi>tt>
SPEC <>3m>em>
C" SUBCAT <elbt.<yB>cm>
CONT [3] <U<B-M>*>
Ef. . ". LAUGHER [2] <apro>
IKDEX <iadex>
GEKD <feadcr>
KUM <<ng>
PER <third>
DTRS <he*d.coxnp-<truc>
HEAD-DTR <phr>*<>
PHOK [7]
SYKSEM <syn>*cm>
LOC <oc> *l>
CAT <c> *t>
MARKING <unm>*rked>
SUBCAT <neli>t.<yn>em>
FIRST [4] <yn>ena>
LOC <loc> *l>
CAT <c> *t>
MARKIKG <<nm>rked>
SUBCAT <di<t-ty>n*>em>
HEAD <BO<B>
SPEC [4]
PRD <boolean>
CASE <nom>
COKT <rei-obj>
RESTR <n>MO>
KAME [1]
NAMED [2]
PARA [2]
REST <eli>*t->yn*>em>
BEAD [6]
CONT [3]
DTRS <he*d-comp.<imc>
HEAD-DTR <U<fb>-w<ord>
PBON <BdUt.pboB.word>
FIRST [4]
REST <clUt.pboB.word>
SYNSEM <tyB*>em>
LOC <loc> *l>
CAT <c> *t>
MARKING <unm>*rked>
SUBCAT <nelwt.<yn>eni>
FIRST [4]
REST <eli>*t.<yn>#<am>
BEAD [5]
CONT [3]
COMP-DTRS <cfi>*t.Mfb>
COMP-DTRS <ndi<u<fb>
FIRST <>*>Bdy-word>
PBON <neli>t.pbon-word>
FIRST [1]
REST <clwt-phoB-word>
SYNSEM [4]
REST <eli>t->if>
<parser> Eater typ< of expr<MioB >ad word* terminated by ';':

```

2.6 Using UTjpc to Display Output

Since HPSG's feature structures quickly grow beyond a manageable size, output functions have been provided that generate MfcX instructions directly from graphs. There axe three macros for

generating L^AT_EX instructions to typeset a feature structure:

(1a-s <graph>) to typeset a feature structure in a miniscule typeface

(1a-m <graph>) to typeset a feature structure in a minute typeface

(1a-l <graph>) to typeset a feature structure in a small typeface

Let's assume that a sentence has just been parsed, and that the command quit has been issued to the parser, placing us back at Lisp top level with the global variable *parse* holding the result of parsing.

```
*** Exiting; first parse found is in *parse*
<parser> Enter type of expression and words terminated by ';':
--> quit;
NIL
> (1a-s *parse*)
??? Enter file name for LaTeX output: ["graphtex.txt"]
```

You are prompted for the name of the file to which the L^AT_EX instructions will be written. Unless you want to prevent it from being overwritten the next time L^AT_EX instructions are generated, just hit return to direct output to the default file:

```
*** Setting LaTeX parameters for style TABULAR in size SMALL
*** Writing LaTeX source to file "graphtex.txt"
*** Done writing LaTeX output
NIL
>
```

Next, type (latex) to start L^AT_EX on the newly written output, and to produce a dvi-file:¹

¹This function is operating system dependent, and might not work in a different implementation of Lisp.

```

> (latex)
*** LaTeXing root file crammedoutput.tex
This is TeX, Version 2.93 (eas «(#) 1.61 88/10/11 17:57:06)
(preloaded format=lplain 88.10.11)
(crammedoutput.tex
LaTeX Version 2.09 <4 Aug 1988>
(/usr/lib/tex/macros/article.sty
Document Style 'article' <16 Mar 88>,
(/usr/lib/tex/macros/artll.sty) (/usr/lib/tex/macros/fleqn.sty)
(/usr/lib/tex/macros/legno.sty) (crammedoutput.aux) (graphtex.txt) [1]
(crammedoutput.aux)
Output written on crammedoutput.dvi (1 page, 5008 bytes).
Transcript written on crammedoutput.log.
**• Starting xdvi on file crammedoutput.dvi
??? Do you want to print the output ?: ["Ho"]
HIL
>

```

The function runs &TeX on the instructions generated earlier, starts up xdvi (an X-Windows based DVI previewer) on the output, and asks whether to print the dvi-file.² Figure 1 shows the final product of this procedure.

2.7 Changing Search Strategy and other Parameters

Table 1 summarizes general system parameters. (See appendix D for StTjX-output related parameters.)

Name	Default Value	Other Possible Values
<code>*search*</code>	DEPTH-TIL-SPLIT	BREADTH DEPTH
<code>*max-no-of -sorts*</code>	140	integer < 250
<code>*max-sentence-length*</code>	20	integer < 40
<code>*circular-features*</code>	("SPEC")	list of strings

Table 1: Some System Parameters

The value of `*search*` controls the search mode of the parser. The two numeric variables control the size of certain internal data structures that are set up during grammar conversion and parsing; setting them to values much larger than necessary will cause Lisp to set up unnecessarily large arrays. It is important to add the names of features that introduce cycles in Feature Structures to the value of `*circular-features*`; otherwise, the parser might get stuck in a loop when trying to expand feature structures (see Section 7.8).

²Again, the latter two operations can not be expected to work on machines outside the LCL without changes to the Lisp code in the file latex-output.lisp.

3 About HPSG

This section constitutes a brief review of some relevant aspects of the HPSG formalism. It is not intended to provide exhaustive coverage; for this, the reader is advised to turn to [Pollard, forthcoming, Pollard and Moshier, 1990].

HPSG uses *sorted feature structures* to model the information contained in linguistic objects. Natural Language grammars are written as constraints on feature structures, and "the linguistic entities which correspond to admissible feature structures constitute the predictions of the theory" ([Pollard and Sag, forthcoming]).

Sorted feature structures are a specific kind of graph structure that can be identified in a number of ways. In this report, as in much of the HPSG literature, I use diagrams of Attribute-Value Matrices (AVMs), and a logical feature structure description language is defined below to represent input to the parser. Inside the parser sorted feature structures are represented in yet another way, which is described in Section 7.5.

3.1 Organization of HPSG Grammars

I have divided HPSG grammars into a number of components:

Signature Four kinds of information are specified in the signature:

- The available sorts, features, and relations
- The subsumption ordering on the sorts.
- Which features are carried by the individual sorts ("feature appropriateness").
- Type restrictions on the values of individual features ("sortal restrictions") and on the arguments of relations

Lexicon The lexicon is the repository of all information carried by lexical sorts.

Principles of Grammar Universally valid well-formedness conditions such as the Head Feature Principle and the Subcategorization principle.

Constituent Structure Principles Well-formedness conditions on phrases related to constituent structure. This includes the Immediate Dominance Principle (IDP), and the Constituent Order Principle (COP).

The IDP and COP are not included with the other principles of grammar because they represent information about phrase structure. One possibility to avoid the combinatorial explosion of disjunctions would be to treat these constraints differently than other disjunctive constraints in the grammar; however, such a scheme has not been implemented here.

How does the above correspond to a system of constraints on models of linguistic objects? From a technical point of view, grammars written in the HPSG formalism consist of two components:

- The signature defines a system of types for feature structures, where each type of feature structure corresponds to a type of linguistic object that has been identified in the empirical domain.
- A system of constraints over this type scheme, including relational constraints, represents the constraints on admissible models.

Given a signature that specifies a type scheme of sorted feature structures, the principles of grammar make up a set of recursive constraint equations over the type scheme. For each sort in the sort hierarchy, the equations specify the constraints that have to be met by well-formed instances of the sort. These constraints are called the sort's *definition*.

3.2 The Signature

As described above, the signature defines a type scheme for sorted feature structures. Formally, it has the following parts:

- SORTS is the set of sorts.
- FEATURES is the set of features.
- \subset is an inheritance (or "subsumption") ordering over SORTS
- Appropriate: SORTS \rightarrow (FEATURES \rightarrow SORTS) is a function that defines which features are appropriate for individual sorts, and type restrictions on the values of those features.

Sorts are similar to record datatypes. Here is an example of an entry that might be found in a sort hierarchy:

sign
PHON: list(*phon-word*)