

FUF: THE UNIVERSAL UNIFIER  
USER MANUAL  
VERSION 5.2

MICHAEL ELHADAD  
DEPARTMENT OF COMPUTER SCIENCE  
BEN GURION UNIVERSITY OF THE NEGEV  
84105 BEER SHEVA, ISRAEL  
ELHADAD@CS.BGU.AC.IL

27 June 1993

Copyright ©1993 Michael Elhadad

### **Abstract**

This document is the user manual for FUF version 5.2, a natural language generator program that uses the technique of unification grammars. The program is composed of two main modules: a unifier and a linearizer. The unifier takes as input a semantic description of the text to be generated and a unification grammar, and produces as output a rich syntactic description of the text. The linearizer interprets this syntactic description and produces an English sentence. This manual includes a detailed presentation of the technique of unification grammars and a reference manual for the current implementation (FUF 5.2). Version 5.2 includes novel techniques in the unification allowing the specification of types and the expression of complete information. It also allows for procedural unification and supports sophisticated forms of control.

# Contents

<b>1</b>	<b>Introduction</b>	<b>6</b>
1.1	How to Read this Manual . . . . .	6
1.2	Function and Content of the Package . . . . .	6
<b>2</b>	<b>Getting Started</b>	<b>7</b>
2.1	Main User Functions . . . . .	7
<b>3</b>	<b>FDs, Unification and Linearization</b>	<b>9</b>
3.1	What is an FD? . . . . .	9
3.2	A Simple Example of Unification . . . . .	9
3.3	Linearization . . . . .	12
<b>4</b>	<b>Writing and Modifying Grammars</b>	<b>14</b>
<b>5</b>	<b>Precise Characterization of FDs</b>	<b>15</b>
5.1	Generalities: Features, Syntax, Paths and Equations . . . . .	15
5.2	FDs as Graphs . . . . .	17
5.3	Functional Descriptions vs. First-order Terms . . . . .	19
5.4	Disjunctions: The ALT and RALT Keywords . . . . .	20
5.5	Optional Features: the OPT Keyword . . . . .	20
5.6	Control of the Ordering: the PATTERN Keyword . . . . .	21
5.7	Explicit Specification of Sub-constituents: the CSET Keyword . . . . .	22
5.7.1	Implicit and Incremental CSET Specification . . . . .	23
5.7.2	Unification of Incremental CSET Specifications . . . . .	24
5.8	The Special Value NONE . . . . .	24
5.9	The Special Value ANY - The Determination Stage . . . . .	25
5.10	The Special Value GIVEN . . . . .	25
5.11	The Special Attribute CAT: General Outline of a Grammar . . . . .	26
<b>6</b>	<b>Modular Organization of Grammars: Def-alt and Def-conj</b>	<b>27</b>
6.1	Modular Definition of FDs . . . . .	27
6.2	Drawing the Map of a Grammar . . . . .	28
<b>7</b>	<b>Defining Input for Regression Testing: The Test Facility</b>	<b>31</b>
<b>8</b>	<b>Using Lists in FDs</b>	<b>33</b>
8.1	Encoding Lists as FDs . . . . .	33
8.2	When to Use Lists: an Example . . . . .	34
8.3	Typical List Traversal in FUF . . . . .	34
8.4	Path Notations for Lists . . . . .	36
<b>9</b>	<b>Types in Unification</b>	<b>37</b>
9.1	Why Types? . . . . .	37
9.1.1	Typed features . . . . .	37
9.1.2	Typed Constituents: The FSET Construct . . . . .	39
9.1.3	Procedural Types . . . . .	40
9.2	Typed Features: define-feature-type . . . . .	41
9.2.1	Type Definition . . . . .	41
9.2.2	The under Family of Constructs . . . . .	41
9.2.3	Drawing Type Hierarchies . . . . .	42

9.3	Typed Constituents: the FSET Construct . . . . .	43
9.4	Procedural Types . . . . .	43
<b>10</b>	<b>EXTERNAL and Unification Macros</b>	<b>45</b>
<b>11</b>	<b>Morphology and Linearization</b>	<b>46</b>
11.1	Lexical Categories are not Unified . . . . .	46
11.2	CATegories Accepted by the Morphology Module . . . . .	46
11.3	Accepted Features for VERB, NOUN, PRONOUN, DET, ORDINAL, CARDINAL and PUNC- TUATION . . . . .	47
11.4	Possible Values for Features NUMBER, PERSON, TENSE, ENDING, BEFORE, AFTER, CASE, GENDER, PERSON, DISTANCE, PRONOUN-TYPE, A-AN, DIGIT and VALUE . . . . .	47
11.5	The Dictionary . . . . .	48
11.6	Linearization and Punctuation . . . . .	48
<b>12</b>	<b>Control in FUF</b>	<b>50</b>
12.1	Complexity . . . . .	51
12.2	Indexing . . . . .	52
12.3	Dependency-directed Backtracking and BK-CLASS . . . . .	53
12.4	Lazy Evaluation and Freeze with wait . . . . .	59
12.4.1	Wait and Constituent Traversal . . . . .	61
12.5	Conditional-Evaluation with Ignore . . . . .	62
<b>13</b>	<b>Tracing and Debugging</b>	<b>63</b>
13.1	What it Means to Debug a FUF Program . . . . .	63
13.2	Checking the Validity of FDs and Grammars . . . . .	64
13.3	Fine Tuning Tracing: Overview of FUF Tracing Functions . . . . .	64
13.4	Identifying Possible Bugs: Trace-bp . . . . .	65
13.5	Levels of Tracing . . . . .	65
13.6	Tracing of Alternatives and Options . . . . .	66
13.7	Local tracing with boundaries . . . . .	67
13.7.1	Special Flags %trace-on% and %trace-off% . . . . .	67
13.7.2	The Special Tracing Flag %break% . . . . .	68
13.8	The trace-enable and trace-disable Family of Functions . . . . .	69
13.9	The :demo directive . . . . .	70
13.10	Tracing of Specific Stages of the Unification . . . . .	71
13.10.1	Trace-determine . . . . .	71
13.10.2	Trace-Category and Hyper-trace-category . . . . .	72
13.10.3	Trace-Cset . . . . .	72
13.10.4	Trace-BK-Class . . . . .	72
13.10.5	Trace-Wait . . . . .	73
13.10.6	Trace-Alts . . . . .	73
13.11	Some Advice on FUF Debugging . . . . .	74
13.11.1	Syntax Errors . . . . .	74
13.11.2	Semantic Errors . . . . .	74
13.11.3	Expression of Negative Constraints . . . . .	75
13.11.4	Control . . . . .	75
<b>14</b>	<b>Manipulation of FDs as Data-structures</b>	<b>77</b>
14.1	FD Accessors . . . . .	77
14.2	FD Relocation . . . . .	78
14.3	FD Normalization . . . . .	79
14.4	Lists of FDs . . . . .	80
<b>15</b>	<b>Reference Manual</b>	<b>81</b>
15.1	Unification functions . . . . .	81
15.1.1	*lexical-categories* . . . . .	81
15.1.2	*u-grammar* . . . . .	81
15.1.3	*cat-attribute* . . . . .	81
15.1.4	u . . . . .	82
15.1.5	u-disjunctions . . . . .	82
15.1.6	uni . . . . .	82

15.1.7	uni-string	83
15.1.8	uni-fd	83
15.1.9	unif	83
15.1.10	u-exhaust	84
15.1.11	u-exhaust-top	84
15.1.12	uni-num	84
15.2	Checking	85
15.2.1	fd-syntax	85
15.2.2	fd-sem	87
15.2.3	fd-p	88
15.2.4	grammar-p	88
15.2.5	get-error-pair	88
15.2.6	normalize-fd	88
15.3	Tracing	89
15.3.1	*all-trace-off*	89
15.3.2	*all-trace-on*	89
15.3.3	*trace-determine*	89
15.3.4	*trace-marker*	89
15.3.5	*top*	89
15.3.6	all-tracing-flags	89
15.3.7	internal-trace-off	90
15.3.8	internal-trace-on	90
15.3.9	trace-disable	90
15.3.10	trace-disable-all	90
15.3.11	trace-disable-match	90
15.3.12	trace-enable	90
15.3.13	trace-enable-all	90
15.3.14	trace-enable-match	91
15.3.15	trace-off	91
15.3.16	trace-on	91
15.3.17	trace-determine	91
15.3.18	trace-bk-class	91
15.3.19	trace-category	91
15.4	Complexity	92
15.4.1	avg-complexity	92
15.4.2	complexity	92
15.5	Manipulation of the Dictionary	93
15.5.1	*dictionary*	93
15.5.2	lexfetch	93
15.5.3	lexstore	93
15.6	Linearization and Morphology	94
15.6.1	call-linearizer	94
15.6.2	gap	94
15.6.3	morphology-help	94
15.7	Manipulation of FDs as Data-structures	94
15.7.1	fd-intersection	94
15.7.2	fd-member	94
15.7.3	fd-to-list	95
15.7.4	gdp	95
15.7.5	gdpp	95
15.7.6	top-gdp	96
15.7.7	top-gdpp	96
15.7.8	alt-gdp	96
15.7.9	list-to-fd	97
15.8	fine tuning of the unifier	97
15.8.1	*any-at-unification*	97
15.8.2	*use-given*	97
15.8.3	*use-any*	97
15.8.4	*keep-cset*	97
15.8.5	*keep-none*	97
15.8.6	*agenda-policy*	98

<b>Appendices</b>	<b>99</b>
<b>A Installation of the Package</b>	<b>100</b>
A.1 Finding the Files . . . . .	100
A.2 Porting to a New Machine . . . . .	100
A.3 Packages . . . . .	101
<b>B Advanced Features</b>	<b>103</b>
B.1 Advanced Uses of Patterns . . . . .	103
B.2 Advanced uses of CSET . . . . .	103
B.3 Long Distance Dependencies and the GAP feature . . . . .	104
B.4 Specifying Complex Constraints: the TEST and CONTROL Keywords . . . . .	104
B.5 Copying vs Conflation: The # notation . . . . .	105
<b>C Non Linguistic Applications of the Unifier: FUF as a Programming Language</b>	<b>106</b>
C.1 Dealing with Lists: The Member/Append Example . . . . .	106
C.2 Representing Lists as FDs . . . . .	107
C.2.1 NIL and variables . . . . .	107
C.2.2 The ~ notation . . . . .	107
C.3 Environment and Variable Names vs. FD and Path . . . . .	108
C.4 Procedures vs. Categories, Arguments vs. Constituents . . . . .	108
C.5 The total FD Includes the Stack of all Computation . . . . .	108
C.6 Analogy with PROLOG programs . . . . .	109
C.7 Use of Set Values in Linguistic Applications . . . . .	109
<b>D Non Standard Features of the Implementation and Restrictions</b>	<b>110</b>
D.1 Typed Features . . . . .	110
D.2 The FSET Special Attribute and Typed Constituents . . . . .	110
D.3 User-defined Types . . . . .	110
D.4 Limits on Disjunction in Input . . . . .	110
D.5 Mergeable Constituents in Patterns . . . . .	111
D.6 Indexing of Alternation . . . . .	111
D.7 Test and Control . . . . .	111
D.8 GIVEN and UNDER . . . . .	111
D.9 EXTERNAL Specifications and Macros . . . . .	112
D.10 User-defined CAT Parameter . . . . .	112
D.11 Resource Limited Processing . . . . .	112
D.12 BK-CLASS Control Specification . . . . .	112
D.13 WAIT Control Specification . . . . .	113
D.14 IGNORE Control Specifications . . . . .	113
<b>E Changes from Version to Version</b>	<b>114</b>
E.1 New Features and Modifications in Version 3 . . . . .	114
E.2 New Features and Modification in Version 4 . . . . .	114
E.3 New Features and Modification in Version 5 . . . . .	114

# List of Figures

5.1	FD as a graph . . . . .	18
5.2	Conflation in an FD graph . . . . .	18
5.3	A grammar for conjunction . . . . .	18
6.1	Map of the SURGE grammar using the def-alt syntax . . . . .	30
12.1	Effect of bk-class . . . . .	58
12.2	Determination of the address of failure . . . . .	58

# Chapter 1

## Introduction

### 1.1 How to Read this Manual

This manual is designed to help you use the FUF package and to describe and explain the technique of unification grammars.

The FUF package is made available to people interested in text generation and/or functional unification. It can be used:

- as a front-end to a text generation system, providing a surface realization component. SURGE, a grammar of English with large syntactic coverage written in FUF, is included for that purpose.
- as an environment for grammar development. People interested in expressing grammatical theories or developing a practical grammar can experiment with the unifier and linearizer.
- as an environment for a study of functional unification. Functional unification is a powerful technique and can be used for non-linguistic or non-grammatical applications.

This manual contains material for people falling in any of these categories. It starts with an introduction to functional unification, its syntax, semantics and terminology. The following chapters deal with the “grammar development” tools: tracing and indexing, a presentation of the morphology component and the dictionary. The next two chapters present the novel features of FUF: typing and control facilities. A chapter is devoted to typing in FUF: type definition, user-defined unification methods and expression of complete information. One chapter is devoted to the flow of control specification (indexing, dependency-directed backtracking and goal freezing). Finally the last chapter is a reference manual to the package. One appendix is devoted to possible non-linguistic applications of the formalism, and compares the formalism with programming languages, in particular with PROLOG.

Note that this manual does **not** describe or document the example grammars provided as examples with the unifier. The sample grammars contain a brief documentation on-line and are accompanied by example inputs. The SURGE grammar is documented in a separate manual.

### 1.2 Function and Content of the Package

FUF implements a natural language surface generator using the theory of unification grammars (cf. bibliography for references). It follows most closely the original FUG formalism introduced in [1]. Its input is a Functional Description (fd) describing the meaning of an utterance and a grammar (also described as an fd). The Syntax of fds is fully described in section 5. The output is an English sentence expressing this meaning according to the grammatical constraints expressed by the grammar.

There are two major stages in this process: unification and linearization.

Unification consists in making the input-fd and the grammar “compatible” in the sense described in [1]. It comes down to enriching the input-fd with directives coming from the grammar and indicating word order, syntactic constructions, number agreement and other features.

The enriched input is then linearized to produce an English sentence. The linearizer includes a morphology module handling all the problems of word formation (s’s, preterits, ...).



## Chapter 2

# Getting Started

Appendix A describes how to install the package on a new machine. Contact your local system administrator to learn how to load the program on your system. You should know how to load the example grammars and corresponding inputs.

## 2.1 Main User Functions

Once the system is loaded, you are ready to run the program. If you are in a hurry to try the system, the user functions are:

```
(UNI INPUT &key GRAMMAR NON-INTERACTIVE (LIMIT 10000))
;; by default the grammar used is *u-grammar*
;;      non-interactive is nil
;;      limit is 10000
;; Complete work : unification + linearization. Outputs a sentence.
;;      If non-interactive is nil, a line of statistics is
;;      also printed.
;;      In any case, stops after limit backtracking points.

(UNI-FD INPUT &key GRAMMAR NON-INTERACTIVE (LIMIT 10000))
;; by default the grammar used is *u-grammar*
;;      non-interactive is nil.
;;      limit is 10000
;; Does only the unification. Outputs the enriched fd. This is the
;; function to use when trying the grammars manipulating lists of gr5.l
;; If non-interactive is nil, a line of statistics is also printed.
;; In any case, stops after limit backtracking points.

CL> (uni ir01)
The boy loves a girl.
CL> (uni-fd ir02)
(# # ...)

(UNIF FD &key (GRAMMAR *u-grammar*))
;; by default the grammar used is *u-grammar*
;; As uni-fd but works even if FD does not contain a CAT feature.
```

If you want to change the grammar, or the input you can edit the files defining it, or the function with the same name.

There are two other useful functions for grammar developers: `fd-p` checks whether a Lisp expression is a syntactically correct Functional Description (FD) to be used as an input. If it is not, helpful error messages are given. `grammar-p` checks whether a grammar is well-formed.

NOTE: use `fd-p` on inputs only and `grammar-p` on grammars only.

```
(FD-P FD &key (PRINT-MESSAGES t) (PRINT-WARNINGS t))
;; --> T if FD is a well-formed FD.
;; --> nil (and error messages) otherwise.
;; The error messages and warnings are only printed if PRINT-MESSAGES and
;; PRINT-WARNINGS are true.
;; DO NOT USE FD-P ON GRAMMARS
```

```

(GRAMMAR-P &optional (GRAMMAR *u-grammar*)
  &key (PRINT-MESSAGES t) (PRINT-WARNINGS t))
;; --> T if GRAMMAR (by default *u-grammar*) is a well-formed grammar.
;; --> nil (and error messages) otherwise.
;; - FD is *u-grammar* by default
;; - PRINT-MESSAGES is t by default.
;; If it is non-nil, some statistics on the grammar are printed.
;; It should be nil when the function is called non-interactively.
;; - PRINT-WARNINGS is nil by default.
;; If it is non-nil, warnings are generated for all paths in the
;; grammar. (It is sometimes a good idea to manually check that all
;; paths are valid.)

```

```

;; Examples:

CL> (fd-p '((a 1) (a 2)))
-----> error, attribute a has 2 incompatible values: 1 and 2.
nil
CL> (grammar-p)
-----> t
CL> (grammar-p '((a 1) (b 2)))
-----> error, a grammar must be a valid FD of the form:
      ((alt (((cat c1)...) ... ((cat cn) ...))). nil.

```

## Chapter 3

# FDs, Unification and Linearization

In this section, we informally introduce the concepts of FDs and unification. The next section provides a complete description of the FDs as used in the package, and presents all available unification mechanisms.

### 3.1 What is an FD?

An FD (functional description) is a data structure representing constraints on an object. It is best viewed as a list of pairs (attribute value). Here is a simple example:

```
((article "the") (noun "cat"))
```

There is a function called `fd-p` in the package that lets you know whether a given Lisp expression is a valid FD or not and gives you helpful error messages if it is not. In FUGs, the same formalism is used for representing both the input expressions and the grammar.

### 3.2 A Simple Example of Unification

We present here a minimal grammar that contains just enough to generate the simplest complete sentences. It is included in file `gr0.1` in the directory containing the examples. A little more complex grammar, handling the active/passive distinction, is available in `gr1.1`, and a more interesting one in `gr2.1`.<sup>1</sup>

```
((alt MAIN (
  ;; a grammar always has the same form: an alternative
  ;; with one branch for each constituent category.

  ;; First branch of the alternative
  ;; Describe the category S.
  ((cat s)
   (prot ((cat np)))
   (goal ((cat np)))
   (verb ((cat vp)
           (number {prot number})))
   (pattern (prot verb goal)))

  ;; Second branch: NP
  ((cat np)
   (n ((cat noun)))
   (alt (
     ;; Proper names don't need an article
     ((proper yes)
      (pattern (n)))
     ;; Common names do
     ((proper no)
      (pattern (det n))
      (det ((cat article)
```

<sup>1</sup>Note that the simplest grammars presented in the manual use the standard phrase structure approach  $S \rightarrow NPVP$ . More advanced grammars use a systemic approach to language (after `gr4`). In general, the FUG formalism is convenient to write systemic grammars, but it can also be used to implement other linguistic models (PS rules, LFG, GPSG or HPSG).

```

      (lex "the" ))))))))
;; Third branch: VP
((cat vp)
 (pattern (v dots))
 (v ((cat verb))))))

```

A few comments on the form of this grammar: the skeleton of a grammar is always the same, a big **alt** (alternation of possible branches, the unifier will pick one compatible branch to unify with the input). Each branch of this alternation corresponds to a single category (here, **S**, **NP** and **VP**).

The second remark is about the form of the input: as shown in the following example, an input is an FD, giving some constraints on certain constituents. The grammar decides what grammatical category corresponds to each constituent.

The next main function of the grammar is to give constraints on the ordering of the words. This is done using the **pattern** special attribute. A **pattern** is followed by a picture of how the constituents of the current FD should be ordered: (**Pattern (prot verb goal)**) means that the **prot** constituent should come just before the **verb** constituent, etc.

In the first branch, the only thing to notice is how the agreement subject/verb is described: the number of the **PROT** will appear in the input as a feature of the FD appearing under **PROT**, as in:

```
(prot ((number plural) (lex "car")))
```

standing for “cars”. To enforce the subject/verb agreement, the grammar picks the feature **number** from the **prot** sub-fd and requests that it be unified with the corresponding feature of the **verb** sub-fd. This is expressed by:

```
(verb ((number prot number)))
```

which means: the value of the **number** feature of **verb** must be the same as the value of the **number** feature of **prot**. The curly-braces notation denotes what is called a “path” which is a pointer within an fd. Note that in this line of the grammar, we refer to **prot number** even though the **prot number** feature does not appear under **prot** in the rest of the grammar. This is a general feature of FUF: any attribute can appear in an FD, and its value can be given either by the grammar directly where it would appear, or by the input, or by the grammar coming from a distant place and using a path.

Note also that the agreement constraint could have been written in the “opposite” direction:

```
(prot ((number verb number)))
```

Or even:

```
(prot number verb number)
```

In the second branch, describing the NPs, we have two cases, corresponding to proper and common nouns. Common nouns are preceded by an article, whereas proper nouns just consist of themselves, *e.g.*, “the car” vs. “John”. If the feature **proper** is not given in the input, the grammar will add it. By default, the current unifier will always try the first branch of an **alt** first. That means that in this grammar, proper nouns are the default.

Finally, a brief word about the general mechanism of the unification: the unifier first unifies the input FD with the grammar. In the following example, this will be the first pass through the grammar. Then, each sub-constituent of the resulting FD that is part of the **cset** (constituent-set) of the FD will be unified again with the whole grammar. This will unify the sub-constituents **prot**, **verb** and **goal** also. This is how recursion is triggered in the grammar. The next section describes how the **cset** is determined. All you need to know at this point is that if a constituent contains a feature (**cat xxx**) it will be tried for unification.

In the input FDs, the sign **==** is used as a shortcut for the notation:

```
(n == John)  <==>  (n ((lex John)))
```

The **lex** feature always contains the single string that is to be used in the English sentence for all “terminal” constituents.

```

;; When unified with the following FD, the grammar will output the
;; sentence ‘‘John likes Mary’’.
(setq ir01 '((cat s)

```

```

      (prot ((n == john)))
      (verb ((v == like)))
      (goal ((n == Mary))))))

;; Which corresponds to the linearization of the following complete
;; FD (this is the result of the unification):

CLISP> (uni-fd ir01)

((cat s)
 (prot ((n ((lex "john")
             (cat noun))))
        (cat np)
        (proper yes)
        (pattern (n))))
 (verb ((v ((lex "like")
             (cat verb))))
        (cat vp)
        (number {prot number})
        (pattern (v dots))))
 (goal ((n ((lex "Mary")
             (cat noun))))
        (cat np)
        (proper yes)
        (pattern (n))))
 (pattern (prot verb goal)))

```

Following the trace of the program will be the easiest way to figure out what is going on:

```

LISP> (uni ir01)
-->
>STARTING CAT S AT LEVEL {}

-->Entering alt TOP -- Jump indexed to branch #1: S matches input S
-->Updating (CAT NIL) with NP at level {PROT CAT}
-->Updating (CAT NIL) with NP at level {GOAL CAT}
-->Updating (CAT NIL) with VP at level {VERB CAT}
-->Enriching input with (NUMBER {PROT NUMBER}) at level {VERB}
-->Enriching input with (PATTERN (PROT VERB GOAL)) at level {}
-->Success with branch #1 S in alt TOP

>STARTING CAT NP AT LEVEL {PROT}

-->Entering alt TOP -- Jump indexed to branch #2: NP matches input NP
-->Updating (CAT NIL) with NOUN at level {PROT N CAT}
-->Enriching input with (NUMBER {PROT NUMBER}) at level {PROT N}
-->Updating (PROPER NIL) with YES at level {PROT PROPER}
-->Enriching input with (PATTERN (N)) at level {PROT}
-->Success with branch #2 NP in alt TOP

>STARTING CAT VP AT LEVEL {VERB}

-->Entering alt TOP -- Jump indexed to branch #3: VP matches input VP
-->Enriching input with (PATTERN (V DOTS)) at level {VERB}
-->Updating (CAT NIL) with VERB at level {VERB V CAT}
-->Success with branch #3 VP in alt TOP

>STARTING CAT NP AT LEVEL {GOAL}

-->Entering alt TOP -- Jump indexed to branch #2: NP matches input NP
-->Updating (CAT NIL) with NOUN at level {GOAL N CAT}
-->Enriching input with (NUMBER {GOAL NUMBER}) at level {GOAL N}
-->Updating (PROPER NIL) with YES at level {GOAL PROPER}
-->Enriching input with (PATTERN (N)) at level {GOAL}
-->Success with branch #2 NP in alt TOP

{Used 3 backtracking points - 0 wrong branches - 0 undos}
John likes mary.

```

In the figure, you can identify each step of the unification: first the top level category is identified: (cat s). The input is unified with the corresponding branch of the grammar (branch #1). Then the constituents are identified. We have here 3 constituents: PROT of cat NP, VERB of cat VP and GOAL of CAT NP. Each constituent is unified in turn. Then for each constituent, the unifier identifies the sub-constituents. In this case,

no constituent has a sub-constituent, and unification succeeds. Note that in general, the tree of constituents is traversed breadth first.

Now, it is also important to know when unification fails. The following example tries to override the subject/verb agreement, causing the failure:

```
(setq ir02 '((cat s)
              (prot ((n == john) (number sing)))
              (verb ((v == like) (number plural)))
              (goal ((n == Mary)))))

LISP> (uni ir02)

>STARTING CAT S AT LEVEL {}

-->Entering alt TOP -- Jump indexed to branch #1: S matches input S
-->Updating (CAT NIL) with NP at level {PROT CAT}
-->Updating (CAT NIL) with NP at level {GOAL CAT}
-->Updating (CAT NIL) with VP at level {VERB CAT}
-->Fail in trying PLURAL with SING at level {VERB NUMBER}

<fail>
```

### 3.3 Linearization

Once the unification has succeeded, the unified fd is sent to the linearizer. The linearizer works by following the directives included in the **pattern**. The exact way to define these features is explained in section 5.6. The linearizer works as follows:

1. Identify the **pattern** feature in the top level: for ir01, it is (**pattern** (**prot verb goal**)).
2. If a pattern is found:
  - (a) For each constituent of the pattern, recursively linearize the constituent. (That means linearize PROT, VERB and GOAL).
  - (b) The linearization of the fd is the concatenation of the linearizations of the constituents in the order prescribed by the pattern feature.
3. If no feature pattern is found:
  - (a) Find the **lex** feature of the fd, and depending on the category of the constituent, the morphological features needed. For example, if fd is of (**cat verb**), the features needed are: **person**, **number**, **tense**.
  - (b) Send the lexical item and the appropriate morphological features to the morphology module. The linearization of the fd is the resulting string. For example, if lex="give" and the features are the default values (as it is in ir01), the result is "gives."

When the fd does not contain a morphological feature, the morphology module provides reasonable defaults. More details on morphology are provided in section 11.

If a pattern contains a reference to a constituent and that constituent does not exist, nothing happens: the linearization of an empty constituent is the empty string. The following example illustrates this feature:

```
;; Unified FD:
((cat s)
 (pattern (prot verb goal benef))
 (prot ((cat noun) (lex "John"))))
 (verb ((cat verb) (lex "like"))))

;; Linearized string (note that constituents GOAL and BENEf are missing):
John likes.
```

Finally, if one of the constituent sent to the morphology is not a known morphological category, the morphology module can not perform the necessary agreements. This is indicated by the following output:

```
;; Unified FD:
((cat s)
 (pattern (prot verb goal))
 (prot ((cat noun) (lex "John"))))
(verb ((cat verb) (lex "like")))
(goal ((cat zozo) (lex "trotteur"))))

;; Linearized string:
John likes <unknown cat ZOZO: trotteur>
```

In general, when you find that in your output, it means that there is something wrong in the grammar. You should check the list of legal morphological categories (see section 11) or you should check why a high level constituent is sent to the morphology (your fd is too flat). You can use the function **morphology-help** to receive on-line help on which categories are known to the morphology module.

## Chapter 4

# Writing and Modifying Grammars

In this section, we briefly outline what steps must be followed to develop a Functional Unification Grammar. The methodology is the following:

1. Determine the input to use. In general, input is given by an underlying application. If not, the criterion to decide what is a good input is that it should be as much “semantic” as possible, and contain the fewest syntactic features as possible.
2. Identify the types of sentences to produce.
3. For each type of sentence, identify the constituents and sub-constituents, and their function in the sentence. A constituent is a group of words that are “tied together” in a clause. A constituent in general plays a certain function with respect to the higher level constituent containing it. For example, in “John gives a book to Mary,” the group “a book” forms a constituent, of category “noun-group,” and it plays the role of the “object upon which action is performed” in the clause. Such role is often called “medium” or “affected” in functional grammars.
4. Determine the output (that is, the unified fds before linearization). In the output, constituents should be grouped in the same pair and the attribute should indicate what function the constituent is fulfilling. In the previous example, we want to have a pair of the form (**medium** [**fd describing** ‘**a book**’]) in the output. The output must also contain all ordering constraints necessary to linearize the sentence and provide all the morphological feature needed to derive all word inflections (*e.g.*, number, person, tense).
5. Determine the “difference” between the input and the output. All features that are in the output but not in the input must be added by the grammar.
6. For each category of constituent, write a branch of the grammar. To do that, you need to specify under which conditions each feature of the “difference” must be added to the input.

This is of course an over-simplified description of the process. Sometimes, the mapping from the input to the output is best considered if decomposed in several stages. For example, in gr4 (cf. file **gr4.1**), the grammar first maps the roles from semantic functions (like **agent** or **medium**) to syntactic roles (like **subject** or **direct-object**), and then does the required syntactic adjustments. In gr11, (cf. file **gr11.1**), there are three stages: first the clause grammar maps from semantic roles to a level called “oblique”, and then oblique is mapped to syntactic functions such as subject or adjunct.

In general, the important idea here is that you must first determine your input and your output and the grammar is the difference of the two.

The process can be complicated if your grammar also includes a lexicon. In this case, a good part of the output should be provided by the lexicon. Grammar gr11 illustrates one way of including the lexicon in your grammar.



## Chapter 5

# Precise Characterization of FDs

### 5.1 Generalities: Features, Syntax, Paths and Equations

An FD is a list of pairs, called features. The attribute of a feature needs to be a symbol. The value of a feature can be either a leaf or recursively an FD. Here is an example:

```
(1) ((cat np)
      (det ((cat article)
            (definite yes))))
      (n ((cat noun)
          (number plural))))
```

A “leaf” is a primitive fd. It can be either a symbol, a number, a string, a character or an array.

A given attribute in an FD must have at most ONE value. Therefore, the FD `((size 1) (size 2))` is illegal. In fact FDs can be viewed as a conjunction of constraints on the description of an object: for an object to be described by `((size 1) (size 2))` it would need to have its property `size` to have both the values 1 and 2. Conversely, if the attribute “size” does not appear in the FD, that means its value is not constrained and it can be anything. The FD `nil` (empty list of pairs) thus represents all the objects in the world. The pair `(att nil)` expresses the constraint that the value of “att” can be anything. It is therefore useless, and the FD `((att1 nil) (att2 val2))` is exactly equivalent to the FD `((att2 val2))`.

Any position in an FD can be unambiguously referred to by the “path” leading from the top-level of the FD to the value considered. For example, FD (1) can be described by the set of expressions:

```
{cat} = np
{det cat} = article
{det definite} = yes
{n cat} = noun
{n number} = plural
```

Paths are represented as simple lists of atoms between curly braces (for example, `{det definite}`). This notation is not ambiguous because at each level there is at most one feature with a given attribute.

A path can be “absolute” or “relative.” An absolute path gives the way from the top-level of the FD down to a value. A relative path starts with the symbol “`^`” (up-arrow). It refers to the FD embedding the current feature. You can have several “`^`” in a row to go up several levels. For example:

```
((cat s)
 (prot ((cat np)
        (number sing)))
 (verb ((cat vp)
        (number {^ ^ prot number}))))

-----|
this is referring to the absolute path {prot number}
```

The notation `{^4 x}` is equivalent to `{^ ^ ^ ^ x}`. It is convenient when dealing with deeply embedded constituents.

Relative paths are not simply a syntactic convenience, but they extend the expressibility of the formalism, by making grammars “relocatable”. For example, the grammar for NPs can be unified with a subconstituent of the input FD at different levels (**agent** and **affected** for example). In each case, a feature like (**determiner** ((**number** {<sup>^</sup> <sup>^</sup> **number**}))) points to the number of the appropriate constituent. Without relative paths such a general constraint could not be expressed.

The value of a pair can be a path. In that case, it means that the values of the pair pointed to by the path and the value of the current pair must always be the same. The two features are then said to be unified. In the previous example, the features at the paths **verb number** and **prot number** are unified. This means that they are absolutely equivalent, they are two names for the same object (structure sharing). This is equivalent to the systemic operation of “conflation”.

In general, an expression of the form  $x = y$ , where either  $x$  or  $y$  is a path or a leaf is called an equation. An fd can be viewed as a flat list of equations.

In FUF, it is possible to have paths on the left of a pair. It is therefore possible to represent an fd as a list of equations as follows:

```
(({ cat } np)
 ({ det cat } article)
 ({ det definite } yes)
 ({ n cat } noun)
 ({ n number } plural))
```

This notation allows to freely mix the “fds as equations” view with the “fds as structure” one.<sup>1</sup>

When using a path on the left, note that the right hand side of the equation is always interpreted as occurring in the context pointed to by the left-hand side. So if you need to use relative paths, the relative path on the right is relative to the end position of the left-hand side. For example, to unify two features {verb syntax number} and {prot number} at level {verb v}, you must write:

```
((verb ((v (({^ syntax number } {^ ^ ^ prot number }))))))
```

and not:

```
((verb ((v (({^ syntax number } {^ prot number }))))))
```

because in this second equation, the path {<sup>^</sup> prot number} is relative to the level {verb syntax number} (not {verb v} as intended) and, therefore, would end up at level {verb syntax prot number} instead of {prot number}.

The only case where a given attribute can appear in several pairs is when it is followed by paths in all but one pairs. That is:

```
((a ((a1 v1)))
 (a {b})
 (a {c}))
```

is a valid FD. It is equivalent for example to:

```
((b ((a1 v1)))
 (a {b})
 (c {b}))
```

or to:

```
((b ((a1 v1)))
 ({a} {b})
 (c {a}))
```

The function **normalize-fd** is convenient to put an FD into its canonical form. For example:

```
(setf fd1 '((a ((a1 v1)))
              (b ((b1 w1)))
              (a ((a2 v2))))
```

<sup>1</sup>Note that the possibility to put paths on the left increases the expressive power of the **external** construct, as it becomes possible to express at run-time constraints on constituents which are not dominated by the position of the external construct in the structure.

```
(b ((b2 ((w2 2))))))
(b ((b2 ((w3 3))))))
```

```
LISP> (normalize fd1)
```

```
((a ((a1 v1)
      (a2 v2)))
 (b ((b1 w1)
      (b2 ((w2 2)
            (w3 3))))))
```

All unification functions assume that the input `fd` is given in canonical form. `normalize-fd` is particularly useful when the inputs are produced incrementally by a program. Note that `normalize-fd` will fail and return `*fail*` if the input FD is not consistent (for example `((a 1) (a 2))`).

## 5.2 FDs as Graphs

When the structure of an FD becomes complex, and more confluents with paths are introduced, a visual representation of the FD becomes extremely useful. This visual representation also provides a clear interpretation of the path mechanism and makes reading of relative path much easier. The structured format of FDs can be viewed as equivalent to a directed graph with labeled arcs as pointed out in [2]. The correspondence is established as follows: an FD is a node, each pair `(attr value)` is a labeled arc leaving this node. The `attr` of the pair is the label of the arc, the value is the adjacent node. Internal nodes in the graph have therefore no label whereas leaves are atomic values. The equivalence is illustrated in Fig.5.1.

The graph notation is particularly useful to interpret relative paths. When a relative path occurs somewhere in an FD, its destination can be identified by going up on the arcs, one arc for each `^`. When the value of a pair is a path, e.g., `(a {b})`, then the corresponding arc actually points to the same node as the given path. In this case, there is structure sharing between `a` and `b`. This configuration is illustrated in Fig.5.2, where the paths `{action number}` and `{agent number}` are conflated, as well as the paths `{affected lex}` and `{affected head lex}` and `{subject}` and `{agent}`.

The conflation of `{subject}` with `{agent}` makes all the paths that are extensions of either `agent` or `subject` equivalent. For example, `{agent lex}` and `{subject lex}` are equivalent. This equivalence is easily read in the graph notation.

The graph notation also makes it clear that the up-arrow notation can be ambiguous. Whenever a Y configuration is met in the graph, for example in the two black nodes in Fig.5.2, the up-arrow does not specify which branch of the Y must be taken. This problem is illustrated in the grammar in Fig.5.3. The FD is extracted from a grammar dealing with conjunction. The constraint enforced by the grammar is that all the conjuncts in a conjunction must have the same syntactic category. A conjunction is represented by an FD with two constituents: `head` represents the conjunction as a whole as a constituent and `list` is a list of conjuncts. The list is represented in a singly-linked list of elements, with a recursive FD containing at each level the first element of the list (feature `car`) and the rest of the list (feature `cdr`). In Fig.5.3, the path `c1` is used to point to the first constituent of the list. `c1` is therefore defined by the equation `c1 = list car`. The grammar in LISP notation is shown below along with a sample input:

```
GR = ((c1 {^ list car})
      (c1 ((cat {^ ^ head cat}))))
```

```
IN = ((head ((cat np)))
      (list ((car ((lex "cat")))
              (cdr ((car ((lex "dog")))
                        (cdr none))))))
```

The expression `(c1 ((cat ...)))` corresponds to the black dot in the graph notation shown in Fig.5.3. The problem is to interpret where the relative path `{^ ^ head cat}` is pointing to. The notation is ambiguous between `{head cat}` and `{list head cat}`, depending on whether one considers the black dot as being located at address `{c1}` or `{list car}`. This ambiguity is solved in FUF by following the convention that up-arrows always refer to the textual location where they appear in the grammar. So in this example, the up-arrows refer to the address `{c1 cat}` and not to the address `{list car cat}` because they are written as a pair `(c1 ((cat {^ ^ head cat})))` and not as `(list ((car ((cat {^ ^ head cat}))))`.

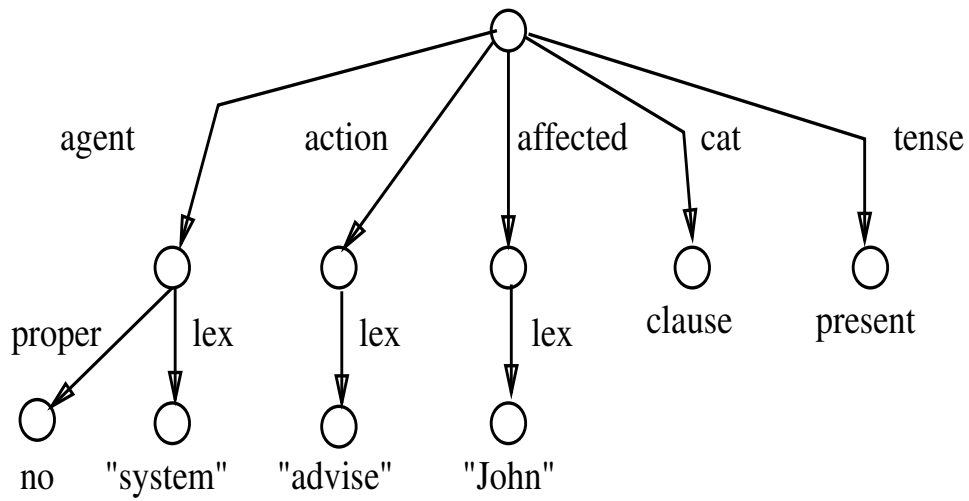


Figure 5.1: FD as a graph

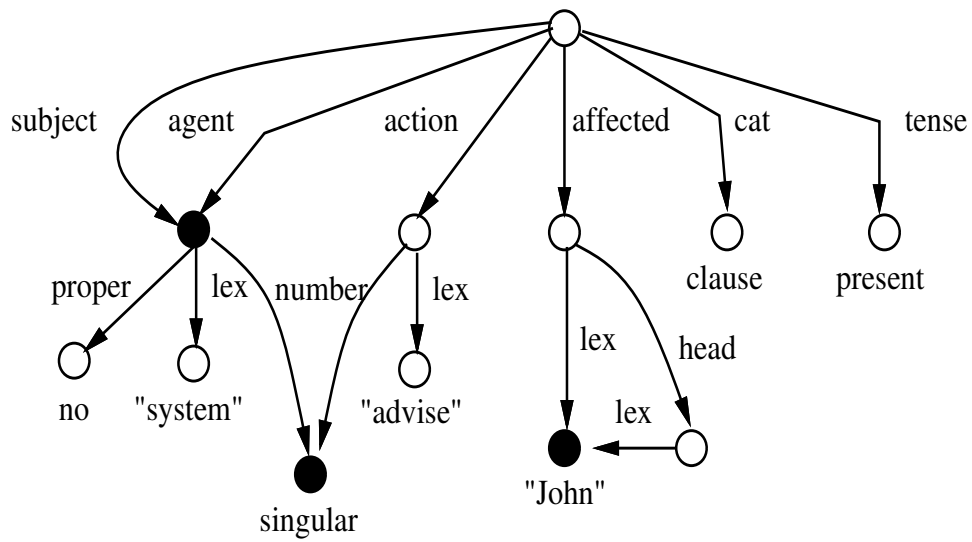


Figure 5.2: Conflation in an FD graph

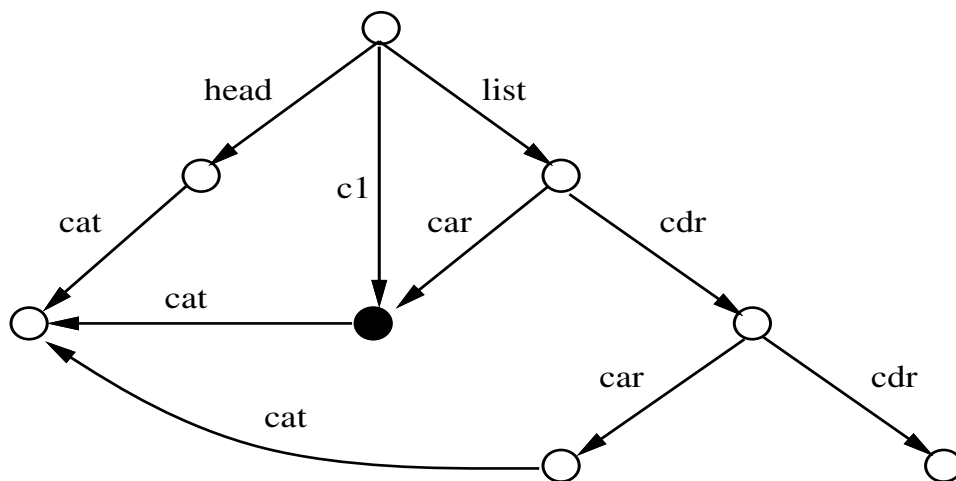


Figure 5.3: A grammar for conjunction

There are special attributes and values which cannot be drawn in this graph notation because they have a special unification behavior. These are, for attributes: **alt**, **opt**, **ralt**, **pattern**, **cset**, **fset**, **test**, **control** and **cat** (or the currently specified **cat** attribute) and for values: **none**, **any** and **given**. The special constructs **#(under x)** and **#(external y)** have also a special meaning for the unifier. These are all the “keywords” known by the unifier. They are presented in the following sections.

### 5.3 Functional Descriptions vs. First-order Terms

To conclude the characterization of FDs as a data-structure, it is useful to contrast functional unification (FU) with the more well known structural unification (SU) as used in PROLOG for example, and to distinguish FDs from the first-order terms used in SU.

The most important difference is that functional unification is not based on order and length. Therefore,  $\{a:1, b:2\}$  and  $\{b:2, a:1\}$  are equivalent in FU but not in SU, and  $\{a:1\}$  and  $\{b:2, a:1\}$  are compatible in FU but not in SU (FDs have no fixed arity). The following quote from Knight summarizes the distinction between feature structures and the first order terms used in SU [3, p.105]:

- Substructures are labeled symbolically, not inferred by argument position.
- Fixed arity is not required.
- The distinction between function and argument is removed.
- Variables and coreference are treated separately.

A comparison between the FD notation and the first-order term notation illustrates these differences. The following FD and first-order term can be used to represent the fact that *Steve builds a crane that is 2 lbs and 4 feet high*:

```
((process build)
 (agent Steve)
 (object ((concept crane)
          (weight 2)
          (height 4))))

build(Steve, Crane(C1, 2, 4))
```

Other representations are of course possible using first-order notation. Some of them have some of the advantages of features structures. For example:

```
build(B1),
agent(B1, Steve),
object(B1, C1),
crane(C1),
weight(C1, 2),
height(C1, 4).
```

In fact, any FD can always be translated in a one-to-one mapping to a class of restricted first-order terms.

Contrasting these two notations for the same example illustrates the differences:

- *Symbolic labels for substructures*: the arguments, that is the **agent** and the **medium**, are clearly labeled in the feature structure notation. In particular, a term like **Crane(C1, 2, 4)** is particularly difficult for a human to interpret.
- *Fixed arity*: features can be added at will to an FD. FDs are used to represent *partial information*. This is not the case for first-order terms. If the knowledge representation changes to include width to crane descriptions, in addition to weight and height, all the terms need to be updated, since **Crane(n, w, h)** is not compatible with **Crane(n, w, h, 1)**. The FD notation is always partial and leaves the possibility of adding new features as needed.
- *Function and argument*: first-order terms have a head (the function) which plays a central role in the unification process. This is not the case in FDs. All information plays the same role.<sup>2</sup>

<sup>2</sup>In most cases, however, one feature plays a special role: for example, the **cat** attribute can specify the category or type of a description, but this is not built into the syntax, and several such “special” attributes can coexist in the same FD, allowing a reader to adopt several perspectives on the same FD.

- *Variables and coreference*: in standard unification, a variable is used to mean two distinct things: that the value of the role is unknown (there are no constraints on it), and that the value of the role is the same as all other objects referred to with the same variable. So for example, in a term such as `like(X, X)`, the use of the variable `X` means that we don't know who likes whom, and that it is known that the agent and the object of `like` must be the same objects. The distinction between these two functions of variables is best explained in [4]. In FDs, coreference and unspecification are represented by two different syntactic devices: variables are features which are unspecified (they simply do not appear in the FD, or appear with an empty value), coreference is handled with *paths*. For example, to express the constraint that the medium of `like` must refer to the same object as its agent, the following FD can be used:

```
((process like)
 (agent {object}))
```

## 5.4 Disjunctions: The ALT and RALT Keywords

`alt` stands for “alternation”. The syntax for using `alt` is:

```
((att1 val1)
 (att2 val2)
 ...
 (ALT {annotations*} (fd1 fd2 ... fdn))
 ...
 (attn valn))
```

The meaning of a pair with an `alt` attribute is the following: the unifier tries to unify the total FD by replacing first the pair `alt` by the FD `fd1`, if this unification fails, then the unifier will try the following alternatives. If all branches of the `alt` fail, the unification fails.

The order in which branches are put within the `alt` does not change the result of the unification. (This is an important feature of the process of unification: the result is always order-independent.) However, since only the first successful unification is returned, order can be used to specify default values. For example, if you want to specify that a sentence should be at the active voice by default, the following order should be used:

```
(ALT (((voice active)
 ... )
 ((voice passive)
 ...)))
```

When the order is truly not relevant and there is no reason to choose a default branch, then you can use the `ralt` keyword instead of `alt`. `ralt` has exactly the same syntax as `alt` and also expresses a disjunction, but the unifier will choose one of the branches at random instead of always trying the first untried branch. (`ralt` stands for “random alt”)

Alternatively, the `:order` annotation can be used to specify whether the branches should be order in random or sequential order. The syntax is as follows:

```
(alt (:order :sequential)      is equivalent to (alt (fd1 ... fdn))
 (fd1 ... fdn))

(alt (:order :random)          is equivalent to (ralt (fd1 ... fdn))
 (fd1 ... fdn))
```

An `alt` can be embedded within another `alt` or it can be the value of a feature as in:

```
((a (alt (1 2 3 4))))
```

## 5.5 Optional Features: the OPT Keyword

`opt` is used to indicate that a set of features is optional. The syntax is

```
((att1 val1)
 ...
 (OPT fd)
 ...
 (attn valn))
```

Its meaning is: if the unification of the whole FD succeeds with `fd`, it is returned as the result. If it fails, the unifier tries again without `fd`. Since the FD `nil` can be unified successfully with any other FD, `opt` is a more readable equivalent to the form:

```
(ALT (fd nil))
```

`opt` is used exactly in the same way as `alt`.

## 5.6 Control of the Ordering: the PATTERN Keyword

As mentioned previously, the generation of a sentence includes two subprocesses: unification and linearization. Unification produces a complex description of a sentence, made of several constituents. Each constituent is described by an FD, and can recursively contain other subconstituents.

Linearization takes such a complex non-ordered description and outputs a linear, ordered, string of words. This operation is constrained by directives put within the FD. These constraints on the ordering appear after the special attribute `pattern`.

For example, in a sentence containing the constituents `prot`, `goal` and `verb`, the following `pattern` can be used:

```
(PATTERN (PROT VERB GOAL))
```

This means that the linearizer should output a string made of the linearization of the constituent `prot` first, followed by the linearization of the constituent `verb` and terminated by the linearization of the constituent `goal`. It also means that nothing can come before `prot` and after `goal`, and nothing can come between each pair.

The constituents correspond to features of the FD describing the sentence. That is, this FD must contain pairs with the attributes `prot`, `verb` and `goal`. For example:

```
(( cat S)
 (PROT (...))
 (GOAL (...))
 (VERB (...))
 (PATTERN (PROT VERB GOAL)))
```

If a constituent mentioned in the pattern is not present in the FD, nothing happens: the linearization of an empty (or non existent) constituent is the empty string.

The `pattern` directives are generally added by the grammar, since the input to the unifier should be a semantic representation and therefore does not contain any constraint on word ordering.

NOTE: Patterns can contain full paths to specify constituents. For example, the following is a legal pattern:

```
(PATTERN ({prot n} {verb v} goal))
```

A given grammar can generate several constraints, that is it can add 2 or more `pattern` pairs to the result. The unifier therefore includes a `pattern` unifier. The role of the pattern unifier is to take several constraints on the ordering and to output one ordering that subsumes all of them.

The following symbols have a special meaning for the pattern unifier: `dots` and `pound` (standing respectively for the notations ‘...’ and ‘#’).

A pattern (`c1 ... c2`) (noted in the program (`c1 dots c2`)) indicates that the constituent `c1` must precede the constituent `c2`, but they need not be adjacent. Zero, one or many other constituents can come in between. The pattern (`c1 ... c2`) still requires the sentence to start with constituent `c1` and to end with `c2`. The pattern (`... c1 ... c2 ...`) only forces `c1` to come before `c2`.

The `pound` (`#`) symbol is used to represent 0 or 1 constituent. For example, if you want to allow a sentence to start with an optional adverbial, you can specify it with the pattern (`# prot ... verb ...`). This directive will be compatible with both (`prot verb goal`) and (`adverb prot verb goal`) for example.

As a consequence of the use of the two symbols `pound` and `dots`, the constraints described by `pattern` directives are PARTIAL orderings.

NOTE: because of the presence of dots and pound, the unification of patterns is a non-deterministic operation. It can produce several results for a given input, and there is no way to predict in which order these possible solutions will be tried. Caution should be exercised when specifying patterns: they should be specific enough to allow only acceptable word orderings (do not use too many **dots**) but should not be too specific to allow for as yet not supported constituents (for example, a sentence can start with an Adverbial, not necessarily an NP).

The following example illustrates the fact that pattern unification is non-deterministic in general:

```
;; Pattern Unification:
p1: (pattern (dots a dots b dots))
p2: (pattern (dots c dots d dots))

;; Compatible Results:
(pattern (dots a dots b dots c dots d dots))
(pattern (dots a dots c dots b dots d dots))
(pattern (dots a dots c dots d dots b dots))
(pattern (dots c dots a dots b dots d dots))
(pattern (dots c dots a dots d dots b dots))
(pattern (dots c dots d dots a dots b dots))

;; Pattern Unification:
p3: (pattern (dots a dots b))
p4: (pattern (dots b c))
;; Pattern Unification fails.
```

Patterns are eventually interpreted by the linearization component to produce a string out of an FD.

Appendix B describes some advanced uses of pattern unification.

## 5.7 Explicit Specification of Sub-constituents: the CSET Keyword

The unifier works top-down recursively: it unifies first the top-level FD against a grammar (generally the top-level FD represents a sentence), and then, recursively, it unifies each of its constituents. For example, to unify a sentence, the unifier first takes the whole FD and unifies it with the grammar of the sentences (**cat S**), then it unifies the **prot** and **goal** with the grammar of NPs (**cat np**), then it unifies the **verb** with the grammar of VPs (**cat vp**).

You can specify explicitly which features of an FD correspond to constituents and therefore need to be recursively unified. To do that, add a pair:

```
(CSET (c1 ... cn))

;; For example:
(CSET (PROT VERB GOAL))
```

The value of a **cset** (stands for Constituent SET) is considered as a SET (unordered). Therefore the following 2 pairs are correctly unified:

```
(CSET (PROT VERB GOAL))
(CSET (VERB GOAL PROT))
```

Actually, two **cset** pairs are unified if and only if there values are two equal sets.

NOTE: A **cset** values can contain full paths to specify constituents. So for example, the following is a legal feature:

```
(cset ({prot n} {verb v} goal))
```

FUF does not rely exclusively on **csets** to find the constituents to be recursively unified. FUF generally tries to infer the value of **cset** from the value of **pattern** and an observation of the features of the current FD (with the assumption that features containing a **cat** attribute are constituents). The exact procedure followed to identify the implicit constituent set of an fd is:

1. If a feature (**cset (c1 ... cn)**) is found in the FD, the constituent set is just (**c1 ... cn**).
2. If no feature **cset** is found, the constituent set is the union of the following sub-fds:



- (a) If a pair contains a feature (`cat xx`), it is considered a constituent.
- (b) If a sub-fd is mentioned in the pattern, it is considered a constituent.

As a consequence, explicit `csets` are rarely necessary. They are generally used when an `fd` contains a sub-`fd` that either is mentioned in the pattern or contains a feature `cat`, but that you do NOT want to unify. In that case, you can explicitly specify the `cset` without including this unwanted sub-`fd`. For larger grammars, however, you should put the emphasis on a clean constituent structure, and therefore you should carefully use the explicit CSET facilities instead of blindly relying on FUF's inferencing. In this case, the advanced CSET facilities described below will prove helpful.

### 5.7.1 Implicit and Incremental CSET Specification

CSET specification is the means by which a programmer describes the constituent structure of sentences. Given the importance of this task and its complexity, facilities have been added to FUF to make CSET specification more flexible. It is thus possible to specify *implicit* and *incremental* CSETs. This section describes these features of the CSET specification. It presents advanced facilities and can be skipped in first lecture.

The general form of a CSET specification is:

```
(CSET (= c1 ... cn)
      (== b1 ... bm)
      (+ a1 ... ap)
      (- d1 ... dq)))
```

The simple syntax `(CSET (c1 ... cn))` is equivalent to `(CSET (= c1 ... cn))`.

As usual in FDs, all of the features of the CSET are optional. Each feature in the CSET feature is named as follows:

- The `=` sublist is called the absolute CSET.
- The `==` sublist is called the basis CSET.
- The `+` and `-` sublists are the increment CSETs.

The idea behind the use of incremental CSETs is to gradually refine the CSET description, by adding partial information - thus folding the constituent structure description into the general “partial information, gradual refinement” methodology of FD specification. The incremental CSET specifications are added to the basis CSET, which can either be specified explicitly, using the `==` notation, or be the result of the implicit CSET inference described above.

The actual CSET of an FD is computed by applying the following procedure:

1. If the absolute CSET (`=`) is present, it becomes the value of the actual CSET. The absolute CSET feature is used when you want to disable all incremental computations.
2. Else: If the basis-CSET is present,
  - (a) let BSET = basis-CSET, else let BSET = the implicit CSET, computed as described above (from pattern and `cat` inspection).
  - (b) If the `+increment`-CSET is present, let BSET = BSET union `+increment`.
  - (c) If the `-increment`-CSET is present, let BSET = BSET - `-increment` (set difference).
  - (d) Return BSET as the actual CSET.

The result of this procedure is a list of absolute paths, pointing to the sub-constituents of the current constituent. In all cases, this actual `cset` is “cleaned up” by applying the following procedure:

- All duplicate paths are removed. Two paths are duplicate if they point to the same FD. That is, even if the paths are distincts but they point to conflated FDs, they are considered duplicates in the CSET.
- All leaf-constituents are removed. A leaf-constituent is a constituent whose value is not a list of pairs - for example, a symbol or a string. The idea is that there is no point in recursing on a leaf-constituent.
- All constituents which are conflated with a `-increment` constituent are also removed from the actual `cset`.

The following grammar fragments illustrate briefly the use of incremental CSET specification.<sup>3</sup>

```
((cat clause)
 (cset (== prot verb))
 (alt (((prot given)
         (subject {^ prot})
         (voice active))
      ((prot none)
         (goal given)
         (cset (- prot) (+ goal)))
      (subject {^ goal}))))))
```

The first cset feature specifies that the basis cset for this clause constituent should *not* be the implicit cset (derived from pattern and cat), but instead the explicit list (prot verb). In the second cset feature, this basis cset is incrementally refined by specifying that when prot is none, prot is not part of the cset and instead goal is added to the cset. If prot is a non-leaf constituent, then the actual cset for this example will be (prot verb), if it is none, then it will be (verb goal).

```
((cat clause)
 (subject ((cat np)))
 (prot {^ subject})
 (cset (- prot)))
```

In this second fragment, no absolute cset is specified (no = feature) and no basis cset is present (no == feature). Therefore the implicit cset is computed and is found to be (subject) (because of the presence of the (cat np) feature under subject). (subject) is thus the initial basis-cset. The incremental cset (- prot) is then added to the cset. The basis-cset remains (subject) at this point. Finally, the cset is “cleaned” - and it is found that prot and subject are duplicates of each other because of the conflation (prot {^ subject}). Since prot is member of the -increment list, its “synonyms” are deleted from the basis cset, and the actual cset is found to be empty ().

The computation of implicit and incremental CSETs can interact with the more advanced control features of goal freezing (with the `wait` construct). These interactions are described in Sect.12.4.1.

### 5.7.2 Unification of Incremental CSET Specifications

When two complete CSET specifications are unified, the following rules are applied: assume that (cset (= a1) (== b1) (+ c1) (- d1)) is unified with (cset (= a2) (== b2) (+ c2) (- d2)),

1. If both a1 and a2 are instantiated:
  - (a) If (set-equal a1 a2) then:
    - i. If c1 union c2 is included in a1, return (cset (= a1)) Else return :fail.
    - ii. If d1 union d2 intersect with a1, return :fail Else return (cset (= a1))
  - (b) else :fail.
2. If only a1 is instantiated and a2 is null:
  - (a) If a1 intersects d1 union d2, then :fail.
  - (b) Else (cset (= a1)) is returned.
3. If neither a1 and a2 is instantiated, compute the unions: b=b1 u b2, c=c1 u c2, d = d1 u d2, and return (cset (== b) (+ c) (- d)).

## 5.8 The Special Value NONE

There is a way to prevent an FD from ever getting a value for a given attribute. The syntax is: (`att NONE`). It means that the FD containing that pair will NEVER have a value for ‘att’. Or in other words, that the object described by the FD has no attribute ‘att’.

<sup>3</sup>Keep in mind that this facility is above all designed for use in large grammars, and the short fragments shown here could easily be handled without incremental CSETs.

## 5.9 The Special Value ANY - The Determination Stage

An **any** value in a pair means that the feature must have a determined value at the end of the unification. A complete unified FD will never contain an **any**, since an **any** stands for something that must be specified. If after unifying everything, the resulting FD contains an **any**, then the unification fails.

An **any** represents a strong constraint. It means that a feature **MUST** be instantiated. **any** should not be understood as “the feature has a value in the input” but as “the feature **WILL** have a value in the result.”

The idea of a “resulting final FD” coming out of the unification is important. It actually implies that the process of unification is the composition of 2 sub-processes: the unification per se and what we call here the “determination.”

The determination process assures that the resulting FD is well formed. It is a necessary stage since the “resulting final” FD is more constrained than regular FDs. Here is what the determination does:

- checks that no **any** is left.
- tests all the **test** constraints.
- tests that no frozen constraint is left.

It is important to realize that none of this can be done before the unification is finished. Section 12.4.1 gives a more complete picture of the determination process and explains why there may be a need for several cycles of determination when `wait` and goal freezing are used.

Note that in practice, ANY is used rarely. The next special value GIVEN is used more often, and is easier to manipulate, except in cases where goal freezing is used.

## 5.10 The Special Value GIVEN

A **given** value in a pair means that the feature must have a real value at the beginning of the unification. A unified fd will never contain a **given** since **given** will always be unified with a real value. **given** is useful to specify what features are necessary in an input. It is also much more efficient than **any**. It is often used in branches of an **alt**, to “test” for the presence of a feature.

The rule is: when you think of using **any**, you often want to use **given** and, conversely, when you use **wait** (for goal freezing) and something does not work, it is because you should use **any** instead of **given**.

**Given** addresses the most obvious limitation of the top-down regime used by FUF when traversing the constituent tree (and computing the cset expansion) in cases when the grammar contains the equivalent of left-recursive rules. **Given** is in a sense the dual of the **any** meta-variable of the original FUG formalism: while **any** requires a feature to be instantiated at the *end* of the unification process, **given** requires it to be instantiated *before* unification starts. Thus, **given** is used to check that a constraint has been specified in the input. A **given** feature can only appear in a grammar, thus, **given** gives a different status to the two arguments of the unification function.<sup>4</sup>

```
(( (cat NP)
  (semr ((possessor GIVEN)))
  (cset (det head))
  (det ((cat NP)
        (possessive yes)
        (semr {^ ^ possessor}))))
...)
```

To illustrate how **given** solves the problem of the left-recursive rules, consider how the possessive NP rule is implemented in FUF in this example. This FD implements the FUF equivalent of a rule:

---

NP  $\rightarrow$  NP( possessive ) head

This is a left-recursive rule which, if used in a top-down control, could lead to infinite recursion of the form  $NP \rightarrow NP(NP(NP....(NPhead...)))$ . Note how the cset expansion in FUF is the equivalent of the rule application in rule-based grammatical formalisms.

---

<sup>4</sup>This means that when **given** is used, unification is no longer symmetric, but it remains order-independent and monotonic.

To avoid using this rule when there is no possessor in the input (which is the step which leads to non-termination), the grammar contains a feature (**possessor GIVEN**), checking that the semantic representation of this constituent does indeed have a possessor specified in the input. If none is specified, then the rule will fail, and the sub-constituent possessive NP will not get created. Without the use of **GIVEN**, this FD would always be successfully unified, and the cset expansion would lead to non-termination. The use of the **given** meta-variable therefore ensures that the top-down regime of FUF is goal-directed.

NOTE: The **under** construct is related to the **given** value. It is presented in Section 9.2.

## 5.11 The Special Attribute CAT: General Outline of a Grammar

Each constituent of an FD is generally characterized by its “category”. In FD terms, this means each constituent includes a feature of the form (**CAT category-name**), where category-name is expected to be an atom.

A grammar is expected to give directives for each possible category, for example NP, VP, or NOUN. The outline of a grammar must be:

```
((alt (
  ((cat s)
    <rest of grammar for category S>)
  ((cat np)
    <rest of grammar for category NP>)
  <other categories>
)))
```

NOTE: The current version of the unifier makes the assumption that the grammar has such a form. The (**CAT xxx**) pairs must appear first. The function **grammar-p** checks that a grammar has the right form. The list of categories known by the grammar can be found by using the function **list-cats**. See appendix D for a list of the non-standard features of this implementation.

NOTE: The symbol identifying categories (**'cat**) can be changed in the program. It is by default **'cat**, but this default can be changed by setting a new value to the **\*cat-attribute\*** variable or by providing an optional argument to the unification functions, as explained in the reference manual.

## Chapter 6

# Modular Organization of Grammars: Def-alt and Def-conj

A large FUG, like many other large computer programs, becomes difficult to read and maintain when its size increases. The size of a FUG is best measured by counting the number of disjunctions it contains, as discussed in Chapter 12.1. The depth of embedding of the disjunctions is also an indication of the complexity. The number of branches in the disjunctive normal form of the grammar accounts for this embedding but is in general not very informative. For example, the grammar of FUF contains 580 disjunctions and has a disjunctive normal size of  $10^{29}$ . Another more conventional measure of the complexity of the grammar viewed as a program, is the number of lines of code: 9,500 lines for a total of 290,000 characters. By all measures, such a grammar is a large program. So the well-known problems of size management in software engineering become an issue for large FUGs: ease of maintenance, readability, modularity etc.

### 6.1 Modular Definition of FDs

A FUG is a large FD containing many alternations. At the top-level, a FUG is conventionally made up of a large alternation where each branch describes a different category. The pattern is:

```
((alt (((cat clause) ...)
      ((cat noun-group) ...)
      ((cat verb-group) ...)
      ...)))
```

Each branch is an embedded FD containing in turn many disjunctions. For example, the top-level clause grammar is made up of the following alternations:

```
((cat clause)
 (alt mood ...)
 (alt transitivity ...)
 (alt voice ...)
 (alt circumstantial ...)
 (alt displaced-constituent ...)
 ...)
```

Each one of these alternation in turn is a large FD, with many embedded alternations.

FUF includes a very simple mechanism to allow the grammar writer to develop such large grammars in a modular way. The key aspect is to allow the grammar writer to abstract away from the details of an FD by giving it a name. Named FDs can then be used anywhere a regular FD is allowed. Actually, it turned out that instead of using named FDs, named features were much more convenient. The new syntax distinguishes between two types of named features: named disjunctions and named conjunctions defined by **def-alt** and **def-conj** respectively. Named features can then be used in any FD as regular features; that is, instead of a pair, an FD can contain a named feature. The syntax is illustrated below:

```
(def-alt number (:index number)
  (((number singular) ...)
   ((number plural) ...)))
```

<pre>(def-conj agreement-3-sing-masc   (number singular)   (person <b>third</b>)   (gender masculine))</pre>	
<i>;; Using named feature</i>	<i>Expanded form</i>
<pre>((alt (((cat noun-group)           (:! number)           ...))))</pre>	<pre>((alt (((cat noun-group)           (alt (:index number)                 ((number singular) ...)                 ((number plural) ...))           ...))))</pre>
<pre>((head ((lex "man")))   (:&amp; agreement-3-sing-masc))</pre>	<pre>((head ((lex "man")))   (number singular)   (person <b>third</b>)   (gender masculine))</pre>

A reference to a named disjunction has the form `(:! name)`; `(:& name)` is used for named conjunction. These forms can appear anywhere a pair could appear in an FD. A named conjunction works like a bundle of features which get spliced in the embedding FD when it is referenced. A named disjunction accepts all the annotations that an alt can carry (control and tracing annotations). The **def-alt** and **def-conj** mechanism works as a macro mechanism for grammars.

This simple syntactic tool allows the use of abstraction in the development of FUGs: by naming parts and hiding levels of details, the structure of the grammar becomes more apparent. Named features can be re-used in different contexts (several places in the grammar and/or in several grammars). Practically, it becomes possible to work separately on a module of the grammar without affecting the other parts; several people can work together on the same grammar; single modules can be re-loaded and re-defined without requiring re-loading the whole grammar. The standard benefits of modularity in programming languages apply to the full extent.

## 6.2 Drawing the Map of a Grammar

Interestingly, the new syntax highlights the similarity between FUGs and the systems of systemic linguistic. FUF includes a tool that draws a graphical map of the grammar by displaying the tree of the named features with their dependents (function **draw-grammar**). A high-level map of the SURGE grammar is shown below. Note how similar this map is to a system in systemic linguistics. The **def-alt** syntax has made this level of organization clearly visible in the grammar without requiring any change to the FUF formalism.

The functions which produce this type of grammar maps in either character mode or in postscript format are described below:

```
DRAW-GRAMMAR (&optional (root *u-grammar*))
;; Draw the map of a grammar, starting at level root (by default the
;; root of the whole grammar).

FUF-POSTSCRIPT (root filename &key (shrink t))
;; Produces a postscript file depicting the map of a grammar.
;; If shrink is t, the map is forced to fit on a single page.

;; Example of output:

> (draw-grammar 'det)
|- DET +
    |- PRE-DET
    |- DEICTIC2
    |- ORDINAL
    |- QUANTIFIER
    |   +
    |   |- QUANT-COUNT-PLURAL
    |   |   +
    |   |   |- QUANT-PARTITIVE
    |   |   |- QUANT-PARTITIVE
    |   |- QUANT-MASS
    |   |   +
    |   |   |- QUANT-PARTITIVE
```

```
|  
|- DET-TYPE +  
|   |- ARTICLE-DET  
|   |- POSSESSIVE-DET  
|   |- QUESTION-DET  
|   |- DEMONSTRATIVE-DET  
|   |- QUANTIFIER-DET
```

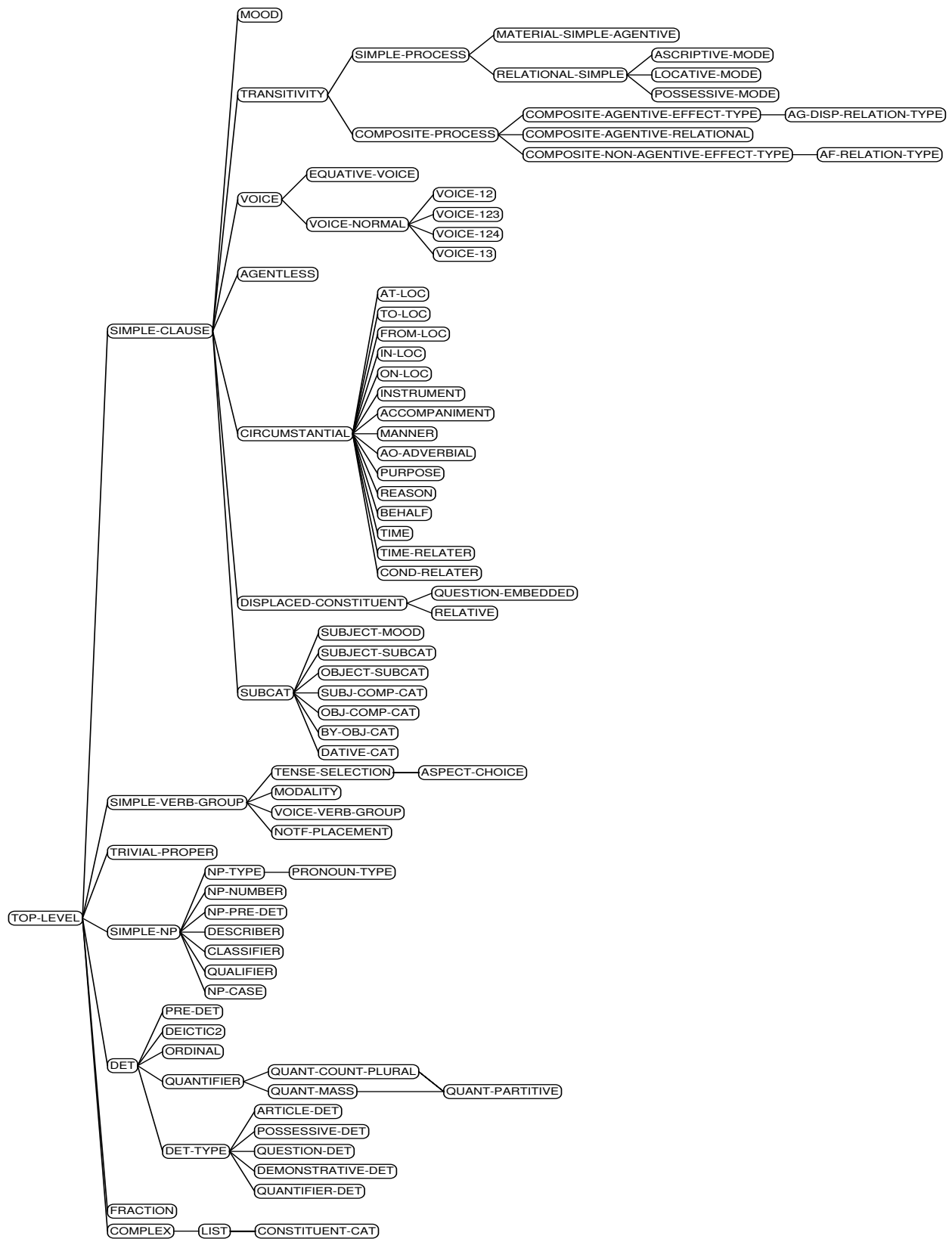


Figure 6.1: Map of the SURGE grammar using the def-alt syntax



## Chapter 7

# Defining Input for Regression Testing: The Test Facility

When developing a grammar, it is important to check that existing input configurations remain compatible whenever the grammar is changed. To this end, FUF includes a regression testing mechanism and a test system. This system allows the grammar developer to define a set of test cases, which are pairs of input FDs together with the sentence they are expected to produce. The test system allows to run a series of tests and checks whether the output is equal to the expected output, it can in addition measure the time spent on each test. A final advantage of the test facility is that test definitions can use the `:&` notation and therefore big FD specifications can be written modularly.

The following functions are defined to take advantage of this facility:

```
(def-test name result input)
;; Define a named test: test on input should produce result.
;; If result is a list, testing result can produce any one of the elements
;; of result.

(get-test name)
;; Return the input for named test.

(clear-tests)
;; Remove all test definitions.

(test &key from to item timed)
;; Evaluate a sequence of tests.
;; A test calls uni-string on input and compares the result with the test's
;; result.
;; If timed is non-nil, time all tests.
;; from and to identify first and last tests in order in which they have
;; been defined.
;; item: identify tests explicitly - either one name or a list of names.

;; Examples:
(def-test :from 't1 :to 't212)
(def-test :item '(t6 t9 t45) :timed t)
(def-test :item 't6)
```

When running the tests, the following output is produced:

```
FUG5 12> (test :item '(t1 t14 t212))

=====
T1 --> "This_car_is_expensive."
{Used 117 backtracking points - 26 wrong branches - 15 undos}
{Used 119 backtracking points - 26 wrong branches - 15 undos}
OK
=====

T14 --> "For_her_to_do_it."
{Used 152 backtracking points - 36 wrong branches - 13 undos}
```

```
{Used 158 backtracking points - 36 wrong branches - 13 undos}
```

```
OK
```

---

```
T212 --> "What."
```

```
{Used 39 backtracking points - 12 wrong branches - 4 undos}
```

```
{Used 40 backtracking points - 12 wrong branches - 4 undos}
```

```
OK
```

---

```
3 tests run - 3 correct.
```

If a test does not succeed, the following output is produced:

```
FUG5 14> (test :item 't1bis)
```

---

```
T1BIS --> "This_car_is_not_expensive."
```

```
{Used 117 backtracking points - 26 wrong branches - 15 undos}
```

```
{Used 119 backtracking points - 26 wrong branches - 15 undos}
```

```
Expected "This_car_is_not_expensive."
```

```
Instead  "This_car_is_expensive."
```

---

```
1 test run - 0 correct.
```

```
The following tests are incorrect: (T1BIS)
```

## Chapter 8

# Using Lists in FDs

Lists of objects are not a primitive type in FUF. The reason is that a list of FDs is not a legal FD. Lists are however very useful in grammatical description, when dealing with subcategorization or conjunction. FUF therefore contains some built-in support for the expression and manipulation of lists. The issue of list manipulation is further developed in Appendix C, where list processing is used as an example showing the expressive power of FUF as a programming language.

This chapter explains how and when to use lists in FUF and describes the FUF facilities easing the use of lists.

### 8.1 Encoding Lists as FDs

A list of FD  $l = (fd_1, fd_2 \dots fd_n)$  is not a legal FD according to the definition of FDs. The main reason why such a form is not accepted in the syntax of FUF is that defining what should be the unification of two lists  $(fd_{11}, fd_{12} \dots fd_{1n})$  and  $(fd_{21}, fd_{22}, \dots fd_{2p})$  is not clear (should it be the union of the 2 lists, its intersection using unification as an equality test, the pairwise unification of the elements one by one, or the product of the unifications, what happens when one unification of the elements fails...) and it is not clear how to extend the notion of path to allow access within a list.

Instead of adding lists as a primitive type of FDs, it is possible to encode lists as regular FDs. This is the approach used and supported in FUF.<sup>1</sup>

Quite simply, lists are represented as FDs with two features, CAR and CDR (with names ala Lisp).

```
;; The list (a b c) is represented by the FD:

((car a)
 (cdr ((car b)
       (cdr ((car c)
             (cdr none))))))

;; The list (a (b c)) is represented by the FD:

((car a)
 (cdr ((car ((car b)
              (cdr ((car c)
                    (cdr none))))))
       (cdr none))))
```

When using this encoding for lists, the path notation can be used to access any element of the list (using a sequence of car and cdr), and the unification method to use when dealing with lists can be defined declaratively in the grammar. FUF includes tools to make the reference to list elements and the development of list-related grammars easier.

---

<sup>1</sup>Note that it is also possible to define lists as procedural types, with any desired specialized unification procedure. But in that case, lists are used as black boxes into which path expressions cannot enter.

## 8.2 When to Use Lists: an Example

Lists should be used whenever you are thinking of using feature names of the form att1, att2 - *i.e.*, whenever you want to go around the limitation of FDs that only one attribute with a given name exist by adding subscripts to the attribute names. Natural candidates for the use of lists are conjunction and the semantic encoding of non-singular objects.

In the following example, lists are used to encode the logical form of a sentence. This encoding encodes in FUF a KL-ONE type of knowledge representation based on entities and binary relations between these entities.

```
;; An FD used in Jacques Robin's Basketball Report Writing System
;; This is a list of entities and a list of relations in which they
;; participate.

;; Semantic content of the sentence:
;; Malone scored 28 points and John Stockton added 27 points and 24 assists.

statistics(malone, 28points)
statistics(stockton, 27points)
statistics(stockton, 24assists)

;; FD notation:

(setf input
  '((cat semantic)

    ;; note the use of the ~ notation to encode a list of FDs.
    (ents ~(((concept player) (name malone))
              ((concept stat) (value 28pts))
              ((concept player) (name stockton))
              ((concept stat) (value 27pts))
              ((concept stat) (value 24asts))))))

    ;; Note the usage of ^n~ to escape to the beginning of the
    ;; containing list and ~n to go down nth elt of a list.
    ;; ^2~ means "go up 2 levels and then to the beginning of the list"
    (rels ~(((concept stat-rel)
                (args ((carrier {^2~ ents ~1})
                      (stat {^2~ ents ~2}))))))
            ((concept stat-rel)
              (args ((carrier {^2~ ents ~3})
                    (stat {^2~ ents ~4}))))))
            ((concept stat-rel)
              (args ((carrier {^2~ ents ~3})
                    (stat {^2~ ents ~4})))))))))
```

The notation `~` and `^n~` are detailed below. They are used to ease the writing of lists as FDs and to traverse lists with path expressions. In this example, note how the lists `ents` (for entities) and `rels` (for relations) are encoded as regular FDs, and how parts of the `ents` list are used in the `rels` FD (for example, the first element of the `rels` list contains two paths to the first and second element of the `ents` list). This form of structure sharing between parts of lists would not be possible using a “black box primitive” list encoding.

## 8.3 Typical List Traversal in FUF

Once lists are encoded as FDs, grammars need to be defined to traverse the lists and unify them according to the specific semantic of each list. Lists are traversed by recursive grammars, using the constituent traversal mechanism implemented by `cset`. The following example shows how to use the input shown above to build a clause out of the semantic input, by composing elements of the appropriate type. The grammar contains the categories `semantic`, `find1` and `find2`. The easiest way to explain how this function works is by taking a procedural interpretation of the FUF flow of control. This procedural interpretation is explained in detail in Appendix C. A constituent can be viewed as a procedure, that receives a certain set of parameters and computes new return values. In this example, the constituent `semantic` receives as input the features `ents` and `rels` which as shown in the FD `input` above. It computes a new feature `clause` as output, which is a linguistic constituent. To this end, it also uses two local variables stored in the features `ent1` and `rel1`. All in all, the function of this category is to map a description of a semantic content into a linguistic clause description.

The grammar also implements two other categories, `find1` and `find2`, which are the ones related to list processing. These categories can be interpreted as procedures searching a list for one or two matches. The simpler

one is `find1`. `Find1` receives two input parameters `in1` and `1st-match` which are a list of FDs and an FD. The category searches the list for the first-match in the list which can be unified with the feature `1st-match`. When it is found, the feature `1st-match` and the element in the list are unified (conflated). The category `find2` operates the same way but for two features `1st-matchA` and `1st-matchB`.

```
;; A grammar to find matching elements in the list of semantic binary
;; relations and compose the matched elements into a linguistic clause
;; structure.
(def-grammar rules ()
  (clear-bk-class)
  '(ALT
    (((cat semantic)
      (ent1 ((cat find2)
        (in2 {^2 ents})
        (1st-matchA ((concept player) (name stockton)))
        (1st-matchB ((concept stat))))))
      (rel1 ((cat find1)
        (in1 {^2 rels})
        (1st-match ((concept stat-rel)
          (args ((carrier {^4 ent1 1st-matchA})
            (stat {^4 ent1 1st-matchB})))))))
      (cset ((= ent1 rel1)))
      (clause ((cat clause)
        (process stat)
        (args {^2 rel1 1st-match args}))))))

;; LIST PROCESSING CONSTITUENTS:
;; Find one element in a list in1 that matches 1st-match and unifies
;; it with 1st-match.
;; If no element in list in1, fail
;; If first element matches 1st-match, unify and stop recursion
;; Else recurse on rest of in1.
;; 1st-match is always linked to the top-level 1st-match feature in
;; the embedded structure of recursive calls.
((cat find1)
  (fset (in1 1st-match rest cat cset))
  (ALT find1 (((1st-match {^ in1 car})
    (cset ()))
    ((rest ((cat find1)
      (1st-match {^2 1st-match})
      (in1 {^2 in1 cdr})))
      (cset ((= rest)))))))

;; Find 2 elements in a list in2 that match 1st-matchA and 1st-matchB
;; and unify them.
((cat find2)
  (fset (in2 1st-matchA 1st-matchB restA restB cat cset))
  (cset ((= restA restB)))
  (restA ((cat find1)
    (1st-match {^2 1st-matchA})
    (in1 {^2 in2})))
  (restB ((cat find1)
    (1st-match {^2 1st-matchB})
    (in1 {^2 in2}))))))

;; To test:
(setq output (uni-fd (input) (rules)))
(top-gdp output {clause args carrier})
(top-gdp output {clause args stat})
```

The list traversal is implemented in the category `find1`. The grammar for `find1` contains a typical list recursion, similar to the following Lisp function:

```
(defun find1 (in1 match1)
  (if (equal match1 (car in1))
      match1
      (find1 (cdr in1) match1)))
```

Naturally, since the grammar is written in FUF it performs differently, but the structure is similar. The `if` test is implemented as an `alt`, and the recursive call is implemented by a `cset` expansion to the sub-constituent `rest`. This structure is typical of the FUF list traversal categories and can be found in a similar form in the SURGE segment dealing with conjunctions.

## 8.4 Path Notations for Lists

To make the use of lists more convenient, three special notations have been defined in FUF:  $\sim$ ,  $\hat{\sim}$  and  $\sim n$ . The macro-character  $\sim$  is used to define lists of FDs as an FD with features `car` and `cdr`. It performs the following transformation:

$$\sim(\text{fd1 fd2} \dots \text{fdn}) \iff ((\text{car fd1})$$
$$\quad (\text{cdr } ((\text{car fd2})$$
$$\quad \quad (\text{cdr} \dots$$
$$\quad \quad \quad ((\text{car fdn})$$
$$\quad \quad \quad (\text{cdr none}))))))$$

The two other notations are used within path expressions (within curly braces).  $\sim n$  is used to access the  $n$ th element within a list. It is expanded to the appropriate sequence of `cdrs` and `car`. The notation  $\hat{\sim}$  is used to go back up to the beginning of a list from within an element of a list. This is expanded at unification-time since it depends on the level of embedding of the element. Therefore  $\hat{\sim}$  actually increases the expressive power of path expressions. If  $\hat{\sim}$  is used within a feature which is not a list element (which means that it does not occur under `car`), then it is equivalent to the simple  $\hat{\sim}$  (which means that it does go up one level anyway). The shorthand notation  $\hat{\sim} n$  is used for  $\hat{\sim} n \hat{\sim}$  that is, go up first  $n$  levels and then to the beginning of the embedding list. These notations are shown for example in the input example above.

# Chapter 9

## Types in Unification

FUF implements three notions of types:

- Typed features
- Procedural types with user-defined unification methods
- Constituent types with the FSET special feature

The idea of using types in unification is relatively recent. So we first motivate the use of types in FUGs. The following sections describe each one of the three methods of typing available.

### 9.1 Why Types?

#### 9.1.1 Typed features

**A Limitation of FUGs: No Structure over the Set of Values** To formally define a grammar, we define  $L$  as a set of labels or attribute names and  $C$  as a set of constants, or simple atomic values. A string of labels (that is an element of  $L^*$ ) is a path, and is noted  $\{l_1 \dots l_n\}$ . A grammar defines a domain of admissible paths,  $D \subset L^*$ .  $D$  defines the skeleton of well-formed FDs.

In functional unification, the set of constants  $C$  has no structure. It is a flat collection of symbols with no relations between each other. All constraints among symbols must be expressed in the grammar. In linguistics, however, grammars assume a rich structure between properties: some groups of features are mutually exclusive; some features are only defined in the context of other features.



Let's consider a fragment of grammar describing noun-phrases (NPs) using the systemic notation given in [5]. Systemic networks, such as this one, encode the choices that need to be made to produce a complex linguistic entity. They indicate how features can be combined or whether features are inconsistent with other combinations. The configuration illustrated by this fragment is typical, and occurs very often in grammars. The schema indicates that a noun can be either a pronoun, a proper noun or a common noun. Note that these three features are mutually exclusive. Note also that the choice between the features **question**, **personal**, **demonstrative**, **quantified** is relevant only when the feature **pronoun** is selected. This system therefore forbids combinations of the type **pronoun**, **proper** and **common**, **personal**.

The traditional technique for expressing these constraints in a FUG is to define a label for each non terminal symbol in the system. The resulting grammar is as follows:

```
((cat noun)
 (alt (((noun pronoun)
 (pronoun
 ((alt (question personal demonstrative quantified))))))
 ((noun proper))
 ((noun common)
 (common ((alt (count mass))))))))
```

This grammar is, however, incorrect, as it allows combinations of the type ((noun proper) (pronoun question)) or even worse ((noun proper) (pronoun zouzou)). Because unification is similar to union of features sets, a feature (pronoun question) in the input would simply get added to the output. In order to enforce the correct constraints, it is therefore necessary to use the meta-FD NONE (which prevents the addition of unwanted features) as shown below:

```
((alt (((noun pronoun)
 (common None)
 (pronoun
 ((alt (question personal demonstrative quantified))))))
 ((noun proper) (pronoun None) (common None))
 ((noun common)
 (pronoun None)
 (common ((alt (count mass))))))))

;; The input FD describing a personal pronoun is then:
((cat noun)
 (noun pronoun)
 (pronoun personal))
```

There are two problems with this corrected FUG implementation. First, both the input FD describing a pronoun and the grammar are redundant and longer than needed. Second, the branches of the alternations in the grammar are interdependent: you need to know in the branch for pronouns that common nouns can be sub-categorized and what the other classes of nouns are. This interdependence prevents any modularity: if a branch is added to an alternation, all other branches need to be modified. It is also an inefficient mechanism as the number of pairs processed during unification is  $O(n^d)$  for a taxonomy of depth  $d$  with an average of  $n$  branches at each level.

**Introducing Typed Features** The problem thus is that FUGs do not gracefully implement mutual exclusion and hierarchical relations. The system of nouns is a typical taxonomic relation. The deeper the taxonomy, the more problems we have expressing it using traditional FUGs.

We propose extracting hierarchical information from the FUG and expressing it as a constraint over the symbols used. The solution is to define a subsumption relation over the set of constants  $C$ . One way to define this order is to define types of symbols, as illustrated below:<sup>1</sup>

```
(define-feature-type noun (pronoun proper common))
(define-feature-type pronoun (personal-pronoun question-pronoun
                             demonstrative-pronoun quantified-pronoun))
(define-feature-type common (count-noun mass-noun))

;; The grammar becomes:
((cat noun)
 (alt (((cat pronoun)
 (cat ((alt (question-pronoun personal-pronoun
             demonstrative-pronoun quantified-pronoun))))))
 ((cat proper))
 ((cat common)
 (cat ((alt (count-noun mass-noun))))))))

;; And the input:
((cat personal-pronoun))
```

The syntax of the new function `define-feature-type` will be presented in section 9.2. Once types and a subsumption relation are defined, the unification algorithm must be modified. The atoms  $X$  and  $Y$  can be unified if they are equal OR if one subsumes the other. The result is the most specific of  $X$  and  $Y$ .

With this new definition of unification, taking advantage of the structure over constants, the grammar and the input become much smaller and more readable. There is no need to introduce artificial labels. The

<sup>1</sup>This notion of typing is similar to the  $\Psi$ -terms defined in [4].



input FD describing a pronoun is a simple `((cat personal-pronoun))` instead of the redundant chain down the hierarchy `((cat noun) (noun pronoun) (pronoun personal))`. Because values can now share the same label `CAT`, mutual exclusion is enforced without adding any pair `1:NONE`.<sup>2</sup> Note that it is now possible to have several pairs `a:vi` in an FD `F`, but that the phrase “the `a` of `F`” is still non-ambiguous: it refers to the most specific of the `vi`. Finally, the fact that there is a taxonomy is explicitly stated in the type definition section whereas it used to be buried in the code of the FUG. This taxonomy is used to document the grammar and to check the validity of input FDs.

### 9.1.2 Typed Constituents: The FSET Construct

A natural extension of the notion of typed features is to type constituents: typing a feature restricts its possible values; typing a constituent restricts the possible features it can have.

```
;; Type declarations (in the grammar):
(determiner ((fset (definite distance demonstrative possessive))))

;; Input FD describing a determiner:
(determiner ((definite yes)
             (distance far)
             (demonstrative no)
             (possessive no)))
```

The `fset` feature specifies that only the four features listed can appear under the constituent `determiner`. This statement declares what the grammar knows about determiners. `Fset` expresses a completeness constraint as defined in LFGs [6]; it says what the grammar needs in order to consider a constituent complete. Without this construct, FDs can only express partial information. The exact syntax of `fset` is given in Section 9.3.

Note that expressing such a constraint (a limit on the arity of a constituent) is impossible in the traditional FU formalism. It would be the equivalent of putting a `NONE` in the attribute field of a pair as in `NONE:NONE`.

```
;; Without FSET:
((cat clause)
 (alt (((process-type action)
         (inherent-roles ((carrier None)
                          (attribute None)
                          (processor None)
                          (phenomenon None)))))
       ((process-type attributive)
        (inherent-roles ((agent None)
                         (medium None)
                         (benef None)
                         (processor None)
                         (phenomenon None)))))
       ((process-type mental)
        (inherent-roles ((agent None)
                         (medium None)
                         (benef None)
                         (carrier None)
                         (attribute None)))))))

;; With FSET:
((cat clause)
 (alt (((process-type action)
         (inherent-roles ((Fset (agent medium benef)))))
       ((process-type attributive)
        (inherent-roles ((Fset (carrier attribute)))))
       ((process-type mental)
        (inherent-roles ((Fset (processor phenomenon)))))))
```

In general, the set of features that are allowed under a certain constituent depends on the value of another feature. Figure 9.3 illustrates the problem. The fragment of grammar shown defines what inherent roles are defined for different types of processes (it follows the classification provided in [7]). We also want to enforce the constraint that the set of inherent roles is “closed”: for an action, the inherent roles are agent, medium and benef *and nothing else*. This constraint cannot be expressed by the standard FUG formalism. `fset` makes it possible.

<sup>2</sup>In this example, the grammar could be a simple flat alternation `((cat ((alt (noun pronoun personal-pronoun ... common mass-noun count-noun))))`, but this expression would hide the structure of the grammar.

Note also that the set of possible features under the constituent **inherent-roles** depends on the value of the feature **process-type**. The first part of the Figure above shows how the constraint can be implemented without **fset**: we need to exclude all the roles that are not defined for the process-type.<sup>3</sup> Note that the problems are very similar to those encountered on the pronoun system: explosion of **none** branches, interdependent branches, long and inefficient grammar.

The **fset** (feature set) attribute solves this problem: **fset** specifies the complete set of legal features at a given level of an FD. **fset** adds constraints on the definition of the domain of admissible paths *D* of a grammar. The syntax is the same as **cset**. Note that all the features specified in **fset** do not need to appear in an FD: only a subset of those can appear. For example, to define the class of middle verbs (*e.g.*, “to shine” which accepts only a medium as inherent role and no agent), the following statement can be unified with the fragment of grammar shown in the previous figure:

```
((verb ((lex "shine")))
 (process-type action)
 (voice-class middle)
 (inherent-roles ((Fset (medium))))))
```

The feature (**FSET (medium)**) can be unified with (**FSET (agent medium benef)**) and the result is (**FSET (medium)**).

Typing constituents is necessary to implement the theoretical claim of LFG that the number of syntactic functions is limited. It also has practical advantages. An important advantage is good documentation of the grammar. Typing also allows checking the validity of inputs as defined by the type declarations.

### 9.1.3 Procedural Types

FUF also implements a third notion of type in unification: procedural types correspond to user-defined data-structures that are unified by special-purpose unification methods. The unification method describes how elements of the type fit in a partial order structure. Typed features are explicitly described (extensionally) partial orders. With procedural types, the partial order is intensionally described by a Lisp procedure.

Procedural types therefore allow the grammar to integrate complex objects that could hardly be described by standard FDs alone. Examples of procedural types are **pattern** (with the pattern unification method enforcing ordering constraints), **cset** (with the cset unification method checking for set equality) and **tpattern** defined in **gr7.1** in the example directory which implements the semantics of tense selection.

There are limitations to the use of procedural types:

1. Procedurally typed objects are always considered as leaves in an FD: that is, no matter how complex is the object, the unifier does not know how to traverse it from the outside. It is viewed as a black box. There is no notion of “path” within the object.
2. Typed objects can only be unified with objects declared of the same type.
3. It is the responsibility of the user to make sure that the unification method actually implements a real partial-order.

The following is a trivial (read: useless) example of how procedural types can be used. The syntax of **define-procedural-type** is described in Section 9.4.

```
;; Unification of 2 numbers is the max: order is the total order of
;; arithmetic (which is also a partial order!)

(defun unify-numbers (n1 n2 &optional path)
  (max n1 n2))

(define-procedural-type 'num 'unify-numbers :syntax 'numberp)

> (u '((num 1)) '((num 2)))
((num 2))

> (u '((num 1)) '((num 0)))
((num 1))

;; Unification of 2 lists is the list with the more elements.
```

<sup>3</sup>Note that this is not even correct, since any other attribute (besides the names of roles) could still be accepted by the grammar.

```
;; That defines a (probably not very useful) total order on lists.
(defun unify-lists (l1 l2 &optional path)
  (if (> (length l1) (length l2))
      l1
      l2))

(define-procedural-type 'list 'unify-lists :syntax 'sequencep)

> (u '((list (1 2 3))) '((list (1 2 3 4))))
((list (1 2 3 4)))

> (u '((list (1 2))) '((list (a b c d))))
((list (a b c d)))
```

## 9.2 Typed Features: define-feature-type

### 9.2.1 Type Definition

Typed features are hierarchies of symbols which are interpreted by the unifier. They are defined by using the function `define-feature-type`.

```
(DEFINE-FEATURE-TYPE <name> <children>)
;; Asserts that <children> are the immediate specializations of <name> in a
;; type hierarchy.

;; Example:

;; finite and non-finite are specializations of the mood symbol.
(define-feature-type mood (finite non-finite))

;; The symbols fact, inference and possible are specializations of the
;; epistemic-modality symbol
(define-feature-type epistemic-modality (fact inference possible))
```

The function `subsume` tests if a symbol (or an object in general) is a specialization of another symbol.

```
(SUBSUME <symbol> <specialization>)
;; T is <specialization> is a specialization of <symbol> NIL otherwise.

;; Example:
(subsume 'mood 'finite)
→ T

(subsume 'epistemic-modality "must")
→ T

(subsume 'finite 'mood)
→ NIL
```

The function `reset-typed-features` resets the working space and deletes all feature type definitions from memory. It is recommended to call it before loading a new grammar to avoid any side effect from previously defined types.

```
> (reset-typed-features)
```

### 9.2.2 The under Family of Constructs

When a type hierarchy is defined, it is possible to check if an input value is more or less “instantiated” within the hierarchy. Using the `under` and `sunder` constructs, one can check if a value is more (resp strictly more) specific than a symbol within a hierarchy. The syntax is the following:

```
((a \#(under <v>)))
;; will unify with
```

```

((a w))
;; only if w is a specialization of v or v itself.

((a \#(sunder <v>)))
;; will unify with
((a w))
;; only if w is a specialization of v but not v itself.

;; The following notations are equivalent:
((a \#(under v))) ((a \#(<= v))) ((a \#(<= v)))

((a \#(sunder v))) ((a \#(< v)))

;; Example:

> (define-feature-type a (aa ab ac))
> (define-feature-type aa (aaa aab))

> (u '((x \#(under aa))) '((x a)))      ;; a is not a specialization of aa
:fail

> (u '((x \#(under aa))) '((x nil)))     ;; nil is the least specific of all
:fail                                     ;; it is not a specialization of aa

> (u '((x \#(<= aa))) '((x aab)))        ;; Ok, aab is a specialization of aa
((x aab))

> (u '((x \#(under z))) '((x z)))        ;; Even if z is not in a hierarchy
((x z))                                 ;; can check for its presence

> (u '((x \#(< z))) '((x z)))            ;; z is not strictly under z.
:fail

```

NOTE: when **under** is used with a symbol **z** which is not part of a type hierarchy, **#(under z)** unifies with **z** only. In particular, it will NOT unify with **NIL**. So the expression **((a #(under z)))** is equivalent to the expression **((a given) (a z))**. In fact, **under** is the typed extension of the notion of **given**. Note that **((a #(under z)))** is the equivalent of LFG's notation  $a =_c z$ .

### 9.2.3 Drawing Type Hierarchies

A pictorial representation of the type hierarchy often helps tremendously to document a grammar. The functions **draw-types** and **types-postscript** create such pictures given the FUF type declarations.

**DRAW-TYPES** (&optional root)  
 Produces a **character** mode picture of the type hierarchy defined in FUF.

**TYPES-POSTSCRIPT** (root filename &key (shrink t))  
 Produces a postscript file depicting the type hierarchy defined in FUF. Use **root = nil** to draw all the types currently defined. Root can also be a **list** of type names, **or** a single type name.

Example of output:

```

> (draw-types 'relation)
|- RELATION +
    |- ASRIPTIVE
    |- POSSESSIVE
    |- LOCATIVE +
        |- SPATIAL
        |- TEMPORAL
        |- ACCOMPANIMENT
        |- EXISTENTIAL
        |- NATURAL-PHENOM

```

>

## 9.3 Typed Constituents: the FSET Construct

The `fset` special attribute expresses a completeness constraint in an FD.

```
;; Value of FSET is a list of symbols.  
(a ((fset (x y z))  
      (x 1))))
```

FSET specifies that all symbols which are NOT listed in its value have a value of `NONE`, that is, they are not defined at this level of the fd.

```
;; b is not in the fset of a  
> (u '((a ((fset (x y z))))) '((a ((b 1)))))  
:fail
```

Two FSET descriptions can be unified. The result is an FSET whose value is the intersection of the two values.

```
> (u '((fset (x y z))) '((fset (v w x z))))  
((fset (x z)))  
  
;; ((fset nil)) is equivalent to NONE (no feature accepted)  
> (u '((a ((fset (x y))))) '((a ((fset (a b)))))  
((a none))
```

## 9.4 Procedural Types

Procedural types are defined by a name and a special unification method, optionally a syntax checker function can be declared. The unification method actually defines a partial order over the elements of the type.

```
(DEFINE-PROCEDURAL-TYPE <name> <function> :syntax <checker> :copier <copier>)  
;; Declares <name> to be a special attribute, whose value can only be  
;; interpreted by <function>.  
  
;; The special types are considered "atomic" types (unifier cannot access to  
;; components from outside).  
  
;; The unification procedure must be \u{deterministic} (no backtracking  
;; allowed) and must be a real "unification" procedure: that is, the type must  
;; be a lattice (or partial order).  
  
;; <FUNCTION> must be a function of 3 args: the vals to unify and the path  
;; where the result is to be located in the total fd.  
;; It must return :fail if unification fails, otherwise, it must return a  
;; valid object of type <type>.  
  
;; NOTE: <FUNCTION> must be such that NIL is always acceptable as an  
;; argument and is always neutral, ie, (<FUNCTION> x nil) = x.  
;; NOTE: <FUNCTION> must be such that (<FUNCTION> x x) = x  
  
;; <CHECKER> must be a function of 1 arg:  
;; It must return either True if the object is a syntactically correct  
;; element of <TYPE>, otherwise, it must return 2 values:  
;; NIL and a string describing the correct syntax of <TYPE>.  
  
;; <COPIER> must be a function of 1 arg:  
;; it must copy an object of <TYPE> that has no cons in common with its  
;; argument. By default, COPY-TREE is used.  
  
;; NOTE: (<COPIER> x) = (<FUNCTION> x nil)
```

The following example shows one use of procedural types:

```
;; Unification of 2 numbers is the max: order is the total order of  
;; arithmetic (which is also a partial order!)
```

```

(defun unify-numbers (n1 n2 &optional path)
  (max n1 n2))

(define-procedural-type 'num 'unify-numbers :syntax 'numberp)

> (u '((num 1)) '((num 2)))
((num 2))

> (u '((num 1)) '((num 0)))
((num 1))

;; Only values of the attribute num can be unified together...
;; a and num are not compatible!
> (u '((num {a})) nil)
:fail

```

## Chapter 10

# EXTERNAL and Unification Macros

The **External** specification allows the grammar writer to produce the constraints of a grammar in a “lazy” way, specifying pieces of the grammar only when needed. **External** also provides a way of developing “macros” in a grammar.

The mechanism is the following: (u x **external**) stops the unification, call the external function specified in the external construct, and uses the value returned to continue as in (u x **value**). External functions expect one argument: the path where the value they return will be used.

The syntax is the following:

```
(( a EXTERNAL))  
  
;; or  
  
(( a \#(EXTERNAL <function >)))
```

In the short form (**external** only), the external function used is the value of the variable **\*default-external-value\***. Otherwise, the name of the function is explicitly specified.

There are two reasons to use an **external** construct:

1. The same portion of the grammar is repeated over and over in different places. Extract this repeating portion, give it a name as a portion, and use the function as a “macro” in the grammar. An example of this sort can be found in file `gr6.l` in the example directory. The macro is called **role-exists**.
2. There are constraints that are better expressed at run-time, when some other parameters, external to the unification process, have been calculated. The **external** construct actually allows a coroutine-like interaction between two processes. This can be used for example to implement a cooperation similar to the one described in the TELEGRAM system [8] between a planner and the unifier. A similar mechanism can be used in the following setting: a unification-based lexical chooser must interact with a knowledge base to decide what lexical items to use. The input given to the lexical chooser only contains pointers to concepts in the knowledge base. When the lexical chooser must make a decision, it needs more information from the knowledge base. The **external** construct allows the lexical chooser to pull information from the knowledge base *only when it needs it*. Therefore, the input does not need to contain all the information that might be needed but only the entry-points to the knowledge base necessary to identify the additional information that may turn out to be relevant under certain conditions.

## Chapter 11

# Morphology and Linearization

The morphology module (partially written by Jay Meyers USC/ISI) makes many assumptions on the form of the incoming functional description. If you want to use it, you must be aware of the following conventions.

### 11.1 Lexical Categories are not Unified

The categories that are handled by the morphology module can be declared to be “lexical categories”. If a category is a lexical category, it is not unified by the unifier, and it is passed unchanged to the morphology module. The assumption here is that the morphology module will do all the reasoning necessary for these categories.

To declare that a category is lexical, you must call the function `register-category-not-unified`. To find out the list of non-unified categories, call the function `categories-not-unified`.

```
CATEGORIES-NOT-UNIFIED (&optional (cat-attribute *cat-attribute*))  
REGISTER-CATEGORY-NOT-UNIFIED (cat &optional (cat-attribute *cat-attribute*))
```

Note that this information depends on the `cat-attribute`, which by default is `cat` (cf. page 26).

### 11.2 CATegories Accepted by the Morphology Module

The following categories only are known by the morphology module. If a category of another type is sent to the morphology, no agreement can be performed. The output in that case is:

```
<Unknown cat CC: LEX>
```

```
MORPH accepts the following values as the value of the attribute CAT:  
ADJ, ADV, CONJ, MODAL, PREP, RELPRO, PUNCTUATION, PHRASE:  
;;      words are sent unmodified.  
NOUN:  
;;      agreement in number is done.  
;;      irregular plural must be put in the list *IRREG-PLURALS*  
;;      in file LEXICON.L  
PRONOUN:  
;;      agreement done on pronoun-type, case, gender, number,  
;;      distance, person.  
;;      irregular pronouns are defined in file LEXICON.L  
VERB:  
;;      agreement is done on number, person, tense and ending.  
;;      irregular verbs must be put in the list *IRREG-VERBS*  
;;      in file LEXICON.L  
DET :  
;;      agreement is done on number, definite and first letter of  
;;      following word for "a"/"an" or feature a-an of following word.  
ORDINAL, CARDINAL:  
;;      string is determined using value and digit to determine whether  
;;      to use digits or letters.
```



The function `morphology-help` will give you this information on-line if you need it.

### 11.3 Accepted Features for VERB, NOUN, PRONOUN, DET, ORDINAL, CARDINAL and PUNCTUATION

VERB:	
ENDING:	{ROOT, INFINITIVE, PAST-PARTICIPLE, PRESENT-PARTICIPLE}
NUMBER:	{SINGULAR, PLURAL}
PERSON:	{FIRST, SECOND, THIRD}
TENSE :	{PRESENT, PAST}
NOUN:	
NUMBER:	{SINGULAR, PLURAL}
FEATURE:	{POSSESSIVE}
A-AN:	{AN, CONSONANT}
PRONOUN:	
PRONOUN-TYPE:	{PERSONAL, DEMONSTRATIVE, QUESTION, QUANTIFIED}
CASE:	{SUBJECTIVE, POSSESSIVE, OBJECTIVE, REFLEXIVE}
GENDER:	{MASCULINE, FEMININE, NEUTER}
PERSON:	{FIRST, SECOND, THIRD}
NUMBER:	{SINGULAR, PLURAL}
DISTANCE:	{NEAR, FAR}
DET :	
NUMBER:	{SINGULAR, PLURAL}
PUNCTUATION:	
BEFORE:	{",", ";", " ", ":", "(", ")", "...}
AFTER :	{",", ";", " ", ":", "(", ")", "...}
ORDINAL, CARDINAL:	
VALUE:	a number (integer or float, positive or negative)
DIGIT:	{YES, NO}

The feature A-AN is used to indicate exceptions to the rule: normally, a noun starting with a consonant is preceded by the indefinite article “a” and if the noun starts with a vowel, it is preceded by “an.” Some nouns start with a consonant but must still be preceded by “an” (for example, “honor” or acronyms “an RST”). In that case, the feature (a-an an) must be added to the corresponding noun.

### 11.4 Possible Values for Features NUMBER, PERSON, TENSE, ENDING, BEFORE, AFTER, CASE, GENDER, PERSON, DISTANCE, PRONOUN-TYPE, A-AN, DIGIT and VALUE

NUMBER:	{SINGULAR, PLURAL}	Default is SINGULAR.
ENDING:	{ROOT, INFINITIVE, PAST-PARTICIPLE, PRESENT-PARTICIPLE}	Default is none.
PERSON:	{FIRST, SECOND, THIRD}	Default is THIRD.
TENSE :	{PRESENT, PAST}	Default is PRESENT.
BEFORE:	{",", ";", " ", ":", "(", ")", "...}	(any punctuation sign) Default is none.
AFTER :	{",", ";", " ", ":", "(", ")", "...}	Default is none.
CASE:	{SUBJECTIVE, OBJECTIVE, POSSESSIVE, REFLEXIVE}	Default is SUBJECTIVE.
GENDER:	{MASCULINE, FEMININE, NEUTER}	

```

        Default is MASCULINE.

PERSON: {FIRST, SECOND, THIRD}
        Default is THIRD.

DISTANCE: {FAR, NEAR}
        Default is NEAR.

PRONOUN-TYPE: {PERSONAL, DEMONSTRATIVE, QUESTION, QUANTIFIED}
        Default is none.

A-AN:    {AN, CONSONANT}
        Default is CONSONANT.

DIGIT:   {YES, NO}
        Default is YES.

VALUE:   a number.

```

## 11.5 The Dictionary

The package includes a dictionary to handle the irregularities of the morphology only: verbs with irregular past forms and nouns with irregular plural only need to be added to the dictionary.

There is no semantic information within this dictionary. In fact, a more sophisticated form of lexicon should have the form of an FD. This dictionary is a part of the morphological module only.

The way to add information to the lexicon is to edit the values of the special variables `*irreg-plurals*` and `*irreg-verbs*`. These variables are defined in the file `LEXICON.L`. After the modification, you need to execute the function `(initialize-lexicon)`. The best way to do that is to edit a copy of the file `LEXICON.L` and to load it back. After loading it, the new lexicon will be ready to use.

The variable `*irreg-plurals*` is a list of pairs of the form (key plural). The default list starts like this:

```

'(("calf" "calves")
 ("child" "children")
 ("clothes" "clothes")
 ("data" "data")
 ...)
```

The variable `*irreg-verbs*` is a list of 5-tuples of the form: (root present-third-person-singular past present-participle past-participle)

The default value starts as follows:

```

'(("become" "becomes" "became" "becoming" "become")
 ("buy" "buys" "bought" "buying" "bought")
 ("come" "comes" "came" "coming" "come")
 ("do" "does" "did" "doing" "done")
 ...)
```

## 11.6 Linearization and Punctuation

The linearizer interprets the **pattern** ordering constraints and assembles the words of the sentence into a linear string. In addition, the linearizer deals with punctuation and capitalization. The general algorithm followed by the linearizer is:

1. If a feature gap is found, the linearization of the FD is the empty string.
2. Else: Identify the **pattern** feature in the current FD. If a pattern is found:
  - (a) For each constituent of the pattern, recursively linearize the constituent.
  - (b) The linearization of the fd is the concatenation of the linearizations of the constituents in the order prescribed by the pattern feature. (Note that during linearization, dots and pounds are ignored in the pattern.)

3. If no feature pattern and a feature lex is found:

- (a) Find the **lex** feature of the fd, and depending on the category of the constituent, the morphological features needed. For example, if fd is of (cat **verb**), the features needed are: **person**, **number**, **tense**.
- (b) Send the lexical item and the appropriate morphological features to the morphology module .
- (c) Identify the feature punctuation. Punctuation can contain three sub-features: before, after and capitalize. According to the value of these features, append or prepend punctuation to the string computed by the morphology module, and capitalize the string if requested. The linearization of the fd is the resulting string.
- (d) If the current cat is not known by the morphology system, the linearization of the constituent is a string of the form “unknown cat X: S”.

4. If no feature pattern and no feature lex is found, the linearization of the current fd is the empty string.

The linearizer also deals with inserting sequences of punctuation signs, insertion of spaces between words and the “liaison” article “an” as in “an interesting case” or “an RPS”. The following rules are implemented:

1. A space is inserted between each pair of words except when:
  - (a) Do not put space after an opening bracket “(“
  - (b) Do not put space before a closing bracket “)”
  - (c) Do not put space before punctuations.
2. When the indefinite singular article (“a”) is followed either by a word that starts with a vowel or by a word whose FD contains the feature (a-an yes), the form “an” is used instead of “a”. For example, if a word is produced from the FD ((lex “RPS”) (a-an yes)), the string “an RPS” will be produced.
3. The first word of the string produced by the linearizer is capitalized.
4. If an FD contains the feature (punctuation ((capitalize yes))), the string produced by the linearizer is capitalized. If the value of capitalize is no, the string is not capitalize, even if it starts the sentence.
5. If an FD contains the feature (punctuation ((before “,”) (after “,”))), the string produced by the linearizer starts and ends with a comma. Any string can be specified in this feature.
6. Leading punctuations are removed from the final string. (There are 6 punctuation signs ,;:!?)
7. A final period (.) is added to sentence if it does not already end with a punctuation. If the mood of the sentence is a specialization of interrogative, then a final question mark (?) is added.
8. Sequences of punctuations are filtered according to the following rules:

(a)

,	,	->	,
.	.	.	.
;	;	;	;
:	:	:	:
!	!	!	!
?	?	?	?

(d)

:	,	:
:	.	:
:	;	:
:	:	:
:	:	:
:	!	:
:	!	:
:	?	:

(b)

.	,	.	,
.	.	.	.
.	;	.	;
.	:	.	:
.	!	.	!
.	?	.	?

(e)

!	,	!	,
!	.	!	.
!	;	!	;
!	:	!	:
!!	:	!!	:
!?	:	!?	:

(c)

;	,	;
;	.	;
;	;	;
;	:	;
;	!	;
;	?	;

(f)

?	,	?	,
?	.	?	.
?	;	?	;
?	:	?	:
?	!	?	!
??	:	??	:

## Chapter 12

# Control in FUF

Pure functional unification can be too slow for practical tasks. FUF offers several control tools that allow the grammar writer to make a grammar more efficient. This section summarizes how control is specified in FUF.

The approach has been to add annotations to the grammars that can be used by the unifier to improve performance. In a sense, this annotation approach is similar to the “optional type declarations” in Common Lisp. An important constraint that has been maintained is that the annotations do not change the semantics of the grammar, but uniquely the order in which the unifier processes it.

All the control annotations are applied to disjunctions. From the measurements made on FUF working on large grammars, it was found that optimization of conjunctions was not necessary, as the average length of conjunctions over the course of a unification is quite small (on the order of 10) and time spent processing them is quite small. In contrast, the overhead associated with dealing with disjunctions is quite high.

The control techniques for disjunctions implemented in FUF are the following:

- indexing
- dependency-directed backtracking
- lazy-evaluation / freeze
- conditional evaluation

In addition to these control mechanisms, a series of “impure” mechanisms ease the integration of unification-based processing in a larger practical system:

- type hierarchies and procedural typing
- external functions
- “given” checking

When all annotations are considered, the syntax of an **alt** or **ralt** constructs is:

```
(alt { traceflag | (trace on traceflag) | (:trace flag) }
      { (index on path) | (:index path)}
      { (demo str) | (:demo str)}
      { (bk-class class) | (:bk-class class) }
      { (:order {:random | :sequential}) }
      { (:wait <list>) }
      { (:ignore-unless <list>) }
      { (:ignore-when <list>) }
      ( fd1 ... fdn ))

;; The annotations (trace, index, bk-class, order, wait and ignore)
;; can appear in any order.
```

This chapter explains and gives examples on how these control features of FUF operate. To better understand why control annotations can make unification faster it is necessary to understand what makes unification slow. So we start by a discussion of how to measure the complexity of a grammar.

## 12.1 Complexity

The complexity of a grammar can be described by the number of possible paths through it, each path corresponding to the choice of one branch for each alternation.

To understand this measure of complexity, it is useful to define the notion of disjunctive normal form . In general, a grammar is made up of embedded disjunctions, that is, of alts within alts. Embedded disjunctions can always be rewritten into top-level disjunctions. This transformation works as indicated below:

<i>;; Embedded form</i>	<i>Normal form</i>
<code>((size ((alt (1 2 3)))))</code>	<code>((alt (((size 1)) ((size 2)) ((size 3)))))</code>
<code>((alt (((cat ((alt (np clause))) (rank group)) ((cat word) (rank word)))))</code>	<code>((alt (((cat np) (rank group)) ((cat clause) (rank group)) ((cat word) (rank word)))))</code>
<code>((size ((alt (1 2))) (weight ((alt (1 2)))))</code>	<code>((alt (((size 1) (weight 1)) ((size 2) (weight 1)) ((size 1) (weight 2)) ((size 2) (weight 2)))))</code>

The number of branches in the top level alt of an FD in normal form is exponential in the depth of the embedding of alts in the original FD. In practice, for large grammars, the normal form can contain up to  $10^{29}$  branches.

In certain implementations, disjunctions are expanded and kept in extensive form instead of being “resolved” as in FUF, by choosing one branch out of the set of branches. In such an approach, the unification of an FD with a disjunction yields a new disjunction:

<code>unify(FD, ((alt (d1...dn)))) = ((alt (unify FD d1)...(unify FD dn)))</code>
---

If one of the branches `di` is not compatible with `FD`, the branch is canceled, by virtue of the equivalences:

<code>((alt (d1 fail))) = d1</code> <code>((alt ())) = fail</code>
---

where `fail` is the symbolic representation of an incompatible FD. The problem with such an approach is that in practice many disjuncts stay “alive” during the unification process, and the result tends to take on sizes close to the normal form expansion of the grammar, which is clearly unmanageable. The FUF implementation does not follow this approach, although one form of the `wait` construct presented in Section 12 allows the grammar writer to selectively keep certain disjunctions unresolved when desired.

The measure of the complexity of a grammar is thus the number of branches the grammar would have in disjunctive normal form. Indexing the grammar actually divides this measure of complexity by a great number.

In FUF, the functions `complexity` and `avg-complexity` compute different measures of the complexity of a grammar.

<pre>(COMPLEXITY &amp;optional grammar with-index) ;; number of branches of grammar in disjunctive normal form. ;; - By default, grammar is *u-grammar* ;; - By default, with-index is T. When it is T, all indexed alts are ;;   considered as one single branch, when it is nil, they are ;;   considered as regular alts.  (AVG-COMPLEXITY &amp;optional grammar with-index rough-avg) ;; "average" number of branches tried when input contains no constraint. ;; - By default, grammar is *u-grammar* ;; - By default, with-index is T. When it is T, all indexed alts are ;;   considered as one single branch, when it is nil, they are ;;   considered as regular alts.</pre>
---

```
;;- By default, rough-avg is nil. When it is nil, the average of an
;; alt is the sum of the complexity of the half first branches. When
;; it is T, the average is half of the sum of the complexity of all
;; branches.
```

Note that these functions do not currently measure the ambiguity of the patterns included in the grammar.

## 12.2 Indexing

The main problem for FUF is to handle non-deterministic constructs in the grammar. The non-deterministic constructs are currently: **alt**, **ralt**, **opt** and **pattern**. Unification of these constructs with an input can produce several results. Whenever the unifier encounters such a construct, it does not know which of the possible results to choose. For example, when unifying an **alt** there is no way to choose a branch out of the many available in the **alt**. The way the program works is to try each of the possibilities one after the other. When the unification later on fails, the program backtracks and tries the next possibility.

This method is actually a blind search through the space of all the descriptions compatible with the grammar. Indexing is a technique used to guide the search in a more efficient way when more knowledge is available. Other heuristics to control search are discussed in Section 12 and include goal freezing (with the **wait** control) and dependency-directed backtracking (with the **bk-class** construct).

FUF allows indexing of **alt** and **ralt** constructs.<sup>1</sup> Indexing tells the unifier how to choose one branch out of the alternation based on the value of the index only, and without considering the other branches ever. The following example illustrates the technique.

```
;; Example taken from gr4

((alt (:trace process) (:index process-type)
  (((process-type actions)
    ...)
   ((process-type mental)
    ...)
   ((process-type attributive)
    ...)
   ((process-type equative)
    ...)))
...)
```

In the example, the `(:index process-type)` declaration indicates that all the branches of the alternation can be distinguished by the value of the **process-type** feature alone. If the input contains a bound feature **process-type**, it is possible to directly choose the corresponding branch of the alternation. If however the input does not correspond such a feature, it has to go through the **alt** in the regular way, with no jumping around.

This is what happens in the tracing messages for each case:

```
;; If input is:
  ((cat clause) (process-type attributive) ...)
;; Trace message is:
  -->Entering alt PROCESS — Jump indexed to branch 3 ATTRIBUTIVE

;; If input does not contain a feature process-type:
  ((cat clause) (prot John) ...)
;; Trace message is:
  -->No value given in input for index PROCESS-TYPE — No jump
  -->Entering alt PROCESS — Branch #1
```

A grammar is always indexed at the top-level by the **cat** feature (or the currently set cat attribute). It makes more sense to index on the features that will be bound in the input or at the moment the **alt** or **ralt** will get tried, but it never hurts to index an **alt**, so it is recommended to index whatever is indexable.

<sup>1</sup>A **opt** construct is actually an **alt** with 2 branches, one being the trivial nil. It would not make sense to index it. A **pattern** construct is ambiguous because patterns like (...a...b...) and (...c...d...) can be combined in many ways. Actually, it is always more efficient to put patterns at the end of the grammar, because much of the ambiguity generated by these patterns would not change the unification anyway, except when the (\* constituent) device is used. In any case, the equivalent of ‘indexing’ a pattern, that is reducing the ambiguity, is to use as few dots as possible in the patterns.

The indexed feature can be at the top level of all the branches, as in the first example for `process-type`, but it can also be at lower levels, like in the following example:

```
;; Example taken from gr4:

((alt verb-trans (:index (verb transitivity-class))
  (((verb ((transitivity-class intransitive)))
    ...)
    ((verb ((transitivity-class transitive)))
    ...)
    ((verb ((transitivity-class bitransitive)))
    ...)
    ((verb ((transitivity-class neuter)))
    ...))
  ...))
```

If the index identifies more than one branches in the disjunction, then the index will serve to prune the disjunction from all the branches that do not match the index, and the search will continue with the remaining branches as usual. For example:

```
((alt (:index cat)
  (((cat clause)
    (mood declarative)
    ...)
    ((cat clause)
    (mood non-finite)
    ...)
    ((cat np)
    ...)))
  ...)
```

When this disjunction is unified with the input `(cat clause)`, two branches are retained by the index matcher (branches 1 and 2). Branch 3 does not match the index and is therefore eliminated. The unifier will then proceed with the alt as if it only contained branches 1 and 2.

Note that if a branch does not contain a bound value for the index, it is as if it contains the value `NIL`, and it will therefore always be retained in the alt after the index matching.

**WARNING:** index does **not** work if several branches are grouped together with an alt embedded under the key attribute. For example, the following input will fail with the following grammar:

```
((cat clause) (mood finite))

((alt (:index cat)
  (((cat ((alt (clause np vp))))
    (cset ((= head))))
    ((cat ((alt (n v adj))))
    (cset (=))))))
```

The unification fails because `clause` does not match the expression `((alt (clause np vp)))` in the index search. The reason for this limitation is that `index` is a speed optimization feature, and making it support `alt` in some rare cases would make all uses of `index` slower. It has therefore been decided to make index more restricted and more efficient.

## 12.3 Dependency-directed Backtracking and BK-CLASS

`BK-class` implements a version of dependency-directed backtracking [9] specialized to the case of FUF. The `bk-class` construct relies on the fact that in FUF, a failure always occurs because there is a conflict between two values for a certain attribute at a certain location in the total FD. In the example illustrating this section, backtracking is necessary because an equation requires that the value of the feature `manner manner-conveyed` be instantiated, but the actual feature is not. The path `manner manner-conveyed` defines the *address of the failure*.<sup>2</sup>

<sup>2</sup>In an FD, each embedded feature can be viewed as an equation between the path leading to the feature in the total FD and the feature value.

The idea is that the location of a failure can be used to identify the only decision points in the backtracking stack that could have caused it. This identification requires additional knowledge that must be declared in the FUG. More precisely, the FUG writer is first allowed to declare certain paths to be of a certain **bk-class**. Then, in the FUG, every choice points that correspond to this **bk-class** must be explicitly declared.

For example, the following statements are used in grammar gr7:

```
(define-bk-class 'manner-conveyed 'manner)
(define-bk-class 'manner 'manner)
(define-bk-class 'lexical-verb '(ao manner))
(define-bk-class 'ao-conveyed 'ao)
(define-bk-class 'ao 'ao)
```

The first one specifies that any path ending with the symbol **manner-conveyed** is of class **manner**. In addition, we tag in the FUG all **alts** that have an influence on the handling of the manner constraint with a declaration (**bk-class manner**) as in:

```
;; In GR7 — CLAUSE branch — ADJUNCTS region
(alt manner (:demo "Is there a manner role?")
  (:bk-class manner)
  (((manner none))
   ;; can it be realized by other means — delay with ANY
   ((manner ((manner-conveyed any))))

   ;; if cannot be realized any other way, resort to an adverb
   ((manner-comp ((cat adv)
    (concept {^ ^ manner concept})
    (lex {^ ^ manner lex}))))
   (manner ((manner-conveyed yes)))
   (pattern (dots manner-comp verb dots))))

  ;; or to a pp
  ((manner-comp ((cat pp)
    (lex {^ ^ manner lex})
    (concept {^ ^ manner concept})
    (opt ((prep ((lex "with"))))))))
  (manner ((manner-conveyed yes)))
  (pattern (dots verb dots manner-comp dots))))
```

This fragment extracted from grammar GR7 illustrates a possible use of the **bk-class** construct. The example is detailed in a paper accompanying the manual (available in file `.../doc/bk-class`). The fragment above implements the following decisions: if there is no **manner** role specified in the input, then nothing needs to be done. If there is a manner role, then the grammar must decide how to realize it. One option is to realize it as either an adverbial adjunct or as a PP adjunct. Another option, shown in the next figure, is to realize the **manner** role by selecting a verb which would carry the manner meaning. The three possibilities are illustrated by the three sentences, where the constituent realizing the manner role is in italics:

1. The Denver Nuggets *narrowly* beat the Boston Celtics 101-99.
2. The Denver Nuggets beat the Boston Celtics *by a slight margin* 101-99.
3. The Denver Nuggets *edged* the Boston Celtics 101-99.

The fragment of the grammar implementing the choice of the verb is shown in the next figure:

```
;; Lexicon for verbs: mapping concept — lex + connotations
((cat lex-verb)
 (alt verbal-lexicon (:demo "What lexical entry can be used for the verb?")
  (:index concept)
  (:bk-class (voice-class transitivity))
  (
   ;; Need to express the result of a game
   ((concept game-result)
    (transitive-class transitive)
    (voice-class non-middle)
    (process-class action))

   (alt (:bk-class (ao manner))
    (
     ;; The verbs edge or nip express the manner
```



```

(({{manner} ((concept narrow)
  (manner-conveyed yes)))
 (lex ((ralt ("edge" "nip")))))

;; The verbs stun or surprise express an evaluation of agent
(({{AO} ((concept rating)
  (carrier {agent})
  (orientation -)
  (ao-conveyed yes)))
 (lex ((ralt ("stun" "surprise")))))

;; Neutral verbs
((lex ((ralt ("beat" "defeat" "down"))))))

;; More verbs for other concepts
....
((concept move)
 (lex ((ralt ("walk" "run"))))))

```

Examples of input that exercise these parts of the grammar are provided in file `ir7`. A simple example is shown in the next figure:

```

(defun t1 ()
  (format t "t1-->_The_Denver_Nuggets_edged_the_Celtics.~%" )
  (setf t1 '((cat clause)
    (process-type action)
    (process ((concept game-result)))
    (agent ((cat proper)
      (lex "The_Denver_Nugget")
      (number plural)))
    (medium ((cat proper)
      (lex "the_Celtic")
      (number plural)))
    (tense past)
    (manner ((concept narrow))))))

```

The overall flow of control through this grammar concerning the realization of the manner role is as follows: The unifier first processes the input at the clause level. When it reaches the region dealing with adjuncts, it checks the **manner** role. There is a **manner** role in input `t1`, so the grammar has a choice between the three modes of realization listed above - verb, adverb or PP. The grammar first tries to “leave a chance to the verb” - more precisely, it leaves a chance to some other unspecified constituent in the grammar to account for the **manner** role. If the verb happens to convey the manner, then nothing else needs to be done. This case is illustrated by the following trace of unification of example `t1` in file `ir7`.

```

;; Trace of example t1

;; Set up the traces
> (load "$fug5/examples/gr7")
> (load "$fug5/examples/ir7")
> (trace-disable-all)
> (trace-enable-match "manner")
> (trace-enable-match "game-result-lex")
> (trace-bk-class t) \index{trace-bk-class (function)}

;; Run example
> (uni (t1))

t1 --> The Denver Nuggets edged the Celtics.
>STARTING CAT CLAUSE AT LEVEL {}

-->Entering alt MANNER -- Branch #1
-->Fail in trying ((CONCEPT NARROW) (LEX "narrowly")) with NONE
  at level {MANNER}
>Special path {MANNER} caught by class (MANNER) after 1 frame

-->Entering alt MANNER -- Branch #2
-->Updating (MANNER-CONVEYED NIL) with ANY at level {MANNER MANNER-CONVEYED}
-->Success with branch 2 in alt MANNER

>Special path {VERB TRANSITIVE-CLASS} caught by class (TRANSITIVITY)
  after 1 frame

```

```

>STARTING CAT PROPER AT LEVEL {SUBJECT}
>STARTING CAT VERB-GROUP AT LEVEL {VERB}
>STARTING CAT PROPER AT LEVEL {OBJECT}
>STARTING CAT LEX-VERB AT LEVEL {VERB LEXICAL-VERB}

;; choose a verb that expresses the manner
-->Entering alt GAME-RESULT-LEX --- Branch #1
-->Fail in trying NONE with RATING at level {AO CONCEPT}
-->Entering alt GAME-RESULT-LEX --- Branch #2
-->Updating (MANNER-CONVEYED ANY) with YES at level {MANNER MANNER-CONVEYED}
-->Updating (LEX NIL) with "nip" at level {VERB LEX}
-->Success with branch 2 in alt GAME-RESULT-LEX

{Used 59 backtracking points - 12 wrong branches - 0 undos}
The Denver Nuggets nipped the Celtics.

```

But if for any reason (like the interaction between manner and argumentation detailed in the bk-class paper and illustrated in examples found in file ir7 or because the input already specifies the verb) the grammar fails to account for the **manner** role in other constituents, the ANY constraint put under the **manner-conveyed** feature will fail at determination time. At this time, we know we failed because of the **manner-conveyed** feature, so if we backtrack, we want to backtrack to the latest choice that involved the **manner** role realization. With the **bk-class** annotations, the unifier will backtrack *directly* to the last choice point of class **manner**, ignoring all intermediate decisions. This case is illustrated by the trace of the unification of example t2 from file ir7. In t2 the verb is specified and set to be “beat” which does not convey the **manner** role. In this case, the unifier needs to backtrack as illustrated:

```

> (uni (t2))
t2 --> The Denver Nuggets narrowly beat the Celtics.
>STARTING CAT CLAUSE AT LEVEL {}

;; First leave a chance to the verb
-->Entering alt MANNER --- Branch #1
-->Fail in trying ((CONCEPT NARROW) (LEX "narrowly")) with NONE
    at level {MANNER}
-->Entering alt MANNER --- Branch #2
-->Updating (MANNER-CONVEYED NIL) with ANY at level {MANNER MANNER-CONVEYED}
-->Success with branch 2 in alt MANNER

>STARTING CAT PROPER AT LEVEL {SUBJECT}
>STARTING CAT VERB-GROUP AT LEVEL {VERB}
>STARTING CAT PROPER AT LEVEL {OBJECT}
>STARTING CAT LEX-VERB AT LEVEL {VERB LEXICAL-VERB}

;; Now choose the verb: it must be beat
-->Entering alt GAME-RESULT-LEX --- Branch #1
-->Fail in trying NONE with RATING at level {AO CONCEPT}
-->Entering alt GAME-RESULT-LEX --- Branch #2
-->Updating (MANNER-CONVEYED ANY) with YES at level {MANNER MANNER-CONVEYED}
-->Fail in trying "beat" with "nip" at level {VERB LEX}
-->Fail in trying "beat" with "edge" at level {VERB LEX}
-->Entering alt GAME-RESULT-LEX --- Branch #3
-->Updating (LEX "beat") with "beat" at level {VERB LEX}
-->Success with branch 3 in alt GAME-RESULT-LEX

;; Problem now: manner is not conveyed! Backtrack...
{Used 64 backtracking points - 17 wrong branches - 1 undo}
>Fail in Determine: found an ANY at level {MANNER MANNER-CONVEYED}
>CURRENT SENTENCE: The Denver Nuggets beat the Celtics.

;; Backtrack because of manner-conveyed
>Special path {MANNER MANNER-CONVEYED} caught by class (AO MANNER)
    after 2 frames

;; Which makes lexical-verb fail:
-->Fail in alt GAME-RESULT-LEX at level {VERB LEXICAL-VERB}

;; Now We skip 23 frames on the stack!
;; We go up DIRECTLY to the choice of manner: express it as an adverb
>Special path {VERB LEXICAL-VERB} caught by class (MANNER) after 23 frames
-->Entering alt MANNER --- Branch #3
-->Updating (CAT NIL) with ADV at level {MANNER-COMP CAT}
-->Enriching input with (CONCEPT {MANNER CONCEPT}) at level {MANNER-COMP}
-->Enriching input with (LEX {MANNER LEX}) at level {MANNER-COMP}
-->Updating (MANNER-CONVEYED NIL) with YES at level {MANNER MANNER-CONVEYED}

```

```

-->Unifying (DOTS START DOTS) with (DOTS MANNER-COMP VERB DOTS)
-->Trying pattern : (DOTS START DOTS MANNER-COMP VERB DOTS)
-->Adding constraints : NIL
-->Success with branch 3 in alt MANNER

;; Redo intermediary decisions
>STARTING CAT PROPER AT LEVEL {SUBJECT}
>STARTING CAT VERB-GROUP AT LEVEL {VERB}
>STARTING CAT PROPER AT LEVEL {OBJECT}
>STARTING CAT LEX-VERB AT LEVEL {VERB LEXICAL-VERB}

;; Choose verb again
-->Entering alt GAME-RESULT-LEX -- Branch #1
-->Entering alt GAME-RESULT-LEX -- Branch #2
-->Entering alt GAME-RESULT-LEX -- Branch #3
-->Updating (LEX "beat") with "beat" at level {VERB LEX}
-->Success with branch 3 in alt GAME-RESULT-LEX

;; This time it works
{Used 106 backtracking points - 58 wrong branches - 137 undos}
The Denver Nuggets narrowly beat the Celtics.

```

Note that in this second example, if the `bk-class` feature is disabled (by typing `(clear-bk-class)`), the unifier does not find an acceptable solution after 100,000 backtracking points. This indicates that `bk-class` is not merely an optimization but is often required to make unification practical.

Figure 12.1 illustrates the effect of using `bk-class`. The unifier therefore uses the knowledge that *only* the verb or the adverb can satisfy the manner constraint in a clause to drastically reduce the search space. But, this knowledge is *locally* expressed at each relevant choice point, retaining the possibility of independently expressing each constraint in the FUG.

In general, the determination of the address of failure is more complex and it is necessary to distinguish between *initial failures* and *derived failures*. An initial failure always occurs at a leaf of the total FD, when trying to unify two incompatible atoms. Failures, however, can also propagate up the structure of the total FD. For example, when unifying `((a ((b 1))))` with `((a ((b 2))))` the original address of failure is the path `a b`. When the unifier backtracks, it also triggers a failure at address `a`, which is not a leaf. This type of failure is called a derived failure. In the implementation of `bk-class`, FUF ignores derived failures and directly backtracks to the first choice point whose `bk-class` matches the last initial failure. Determining whether a failure is derived or initial can be difficult. In FUF, the following technique is used: I distinguish between *forward* and *backward* processing. Forward processing is the normal advance of the search procedure through the disjunctions in the grammar. Upon failure, backward processing starts, which means that the stack of the last decisions is unwound. Each time a new disjunction is retried and restarted, forward processing could restart, but unification may fail again immediately, without making any changes to the total FD being processed. In general, backward mode continues until one of two events occurs: either a restarted alt succeeds and at least one change is made to the total FD (a feature is added to the input), or a failure occurs when unifying a grammar feature with a feature that was present in the original input, before unification started. The distinction between forward and backward processing and the way the address of failure is maintained is illustrated in Fig.12.2.

The figure shows the search space when traversing two successive alts in the grammar. Each arrow corresponds to either the choice of a branch in the alt or to a forced backtracking. Arrows 1 and 2 correspond to the first forward processing step and end up in a failure. This is an initial failure and the *address of failure* (which is by definition the location which caused the latest significant failure) is set to the current position in the total FD. The backward processing step starts and a second branch of the ALT is tried (arcs 3 and 4). This second branch fails immediately, before any change is made to the total FD. In such a case, the backward processing step continues, and the address of failure is **not** updated. The motivation is that the unifier is still backtracking for the same original reason: nothing changed in the grammar that makes the original failure less relevant. So the unifier keeps backtracking, and tries the third and last branch of the ALT (arcs 5 and 6). The same failure occurs again, and backward processing proceeds. The second ALT is exhausted and, therefore, fails, and the second branch of the first ALT is tried in turn (arcs 7, 8 and 9). At this time, two reasons can trigger a switch of failure address:

- Unification fails immediately, before any changes are made, but the failure is due to the incompatibility of a feature of the grammar with a feature that was originally present in the input before unification started. In such a case, the failure is not “caused” by the branch of processing that led to the initial failure, it is caused by the grammar itself with the original input, and would have had occurred even if no decisions had been taken before. So the failure address is switched to the current address.

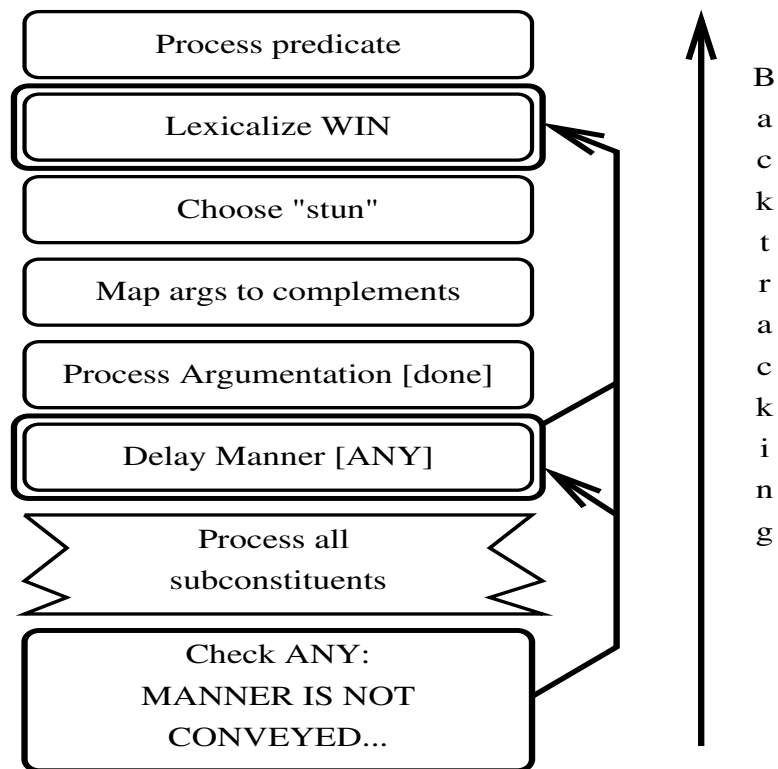


Figure 12.1: Effect of bk-class

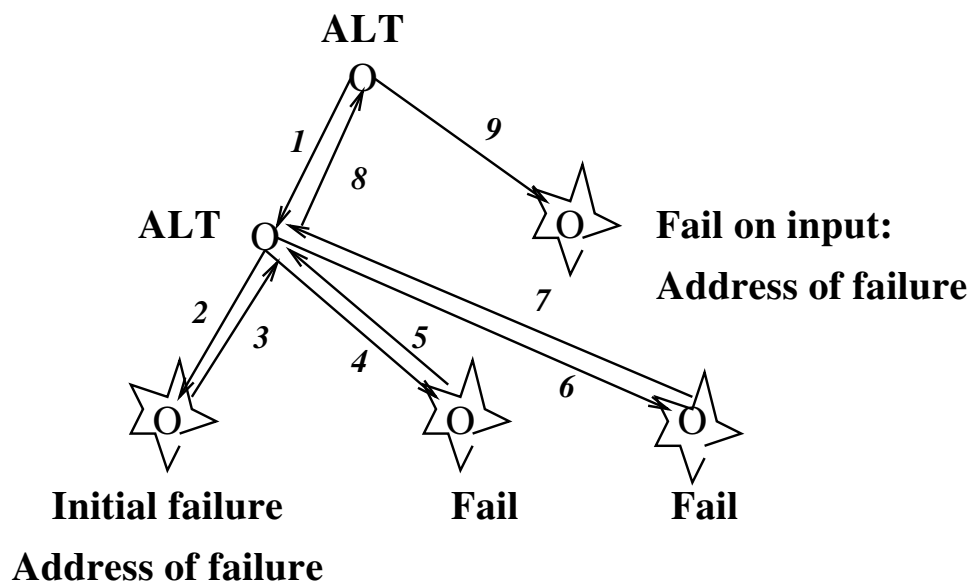


Figure 12.2: Determination of the address of failure

- Unification succeeds and at least one feature from the grammar is added to the total FD. Forward processing starts again. Any future failure can be caused by these new changes, and cannot therefore be attributed to the old initial failure. So the failure address is also switched to the current address.

There is one more complication with the use of the address of failure as a source of information on what decision caused the failure: as already mentioned in Sect.5.2, a given location in the total FD can be accessed through several paths from the top of the total FD (since an FD is not a tree but can be an arbitrary graph). So several distinct paths, with different labels can identify the location where the failure occurred. The intelligent backtracking component must decide which path should be matched against the bk-class specifications. The convention used here is the same one used to disambiguate the up-arrow notation presented in Sect.5.2: the path used to identify the address of failure is the path written in the grammar on the feature that caused the failure. The motivation is that the grammar writer uses the path that is the most meaningful to access features, and that incidental confluences along this path would not be as informative for backtracking purposes. Another approach could have been to use all the path addresses that identify a location and match them against all bk-class specifications. But this approach would have imposed too much overhead on the regular backtracking mechanism.

The bk-class mechanism works well in practice. For the basketball example discussed above, I have measured the number of backtracking points required to generate different clauses conveying the same core content. The following table summarizes these measurements.<sup>3</sup>

Input / Output	w/o bk-class	w/ bk-class
No floating constraints: <i>The DN beat the BC</i>	110	110
Manner in the verb: <i>The DN edged the BC</i>	110	110
Manner as adverbial: <i>The DN narrowly beat the BC</i>	over 10,000	214
AO in the verb: <i>The DN stunned the BC</i>	112	112
AO as adjective: <i>The hapless DN beat the BC</i>	1,623	239
AO as adverbial: <i>The DN surprisingly beat the BC</i>	over 100,000	268
AO & manner together: <i>The hapless DN edged the BC</i>	1,178	246

The number of backtracking points required to generate each example clause is listed with and without **bk-class**. The numbers for the first clause, which does not include any floating constraints, give an indication of the size of the grammar. It roughly corresponds to the number of unretracted decisions made by the grammar. It is the optimal number of backtracking points that a search control regime can obtain for the given input with this grammar. Without **bk-class**, the wide variation in number of backtracking points among the examples indicates the exponential nature of the blind search which floating constraints impose on the standard control regime. In contrast, with **bk-class**, the variation in number of backtracking points remains within a factor of three among all the examples.

## 12.4 Lazy Evaluation and Freeze with wait

The **bk-class** mechanism is useful in general to correct directly a decision made a long time ago, as soon as it appears that the decision was wrong. But this implies that all the work done between the wrong decision and the realization that it is wrong must be re-done once the original wrong has been corrected. For example, in the following case:

```

Frame TOP -----> failure because of feature P
Frame T100 )
Frame T099 )
... )-----> Lots of decisions that do not
... )          affect the value of P.
... )
Frame T016 )
Frame T015 -----> feature P is set to wrong value
Frame T014
...
Frame T001

```

The decisions corresponding to frames T016 to T100 are going to be done twice because T015 made the wrong choice originally. While making a decision twice is not too bad compared to the cost of exhaustively searching all the choices corresponding to frames T016 to T100 backwards, it is still sub-optimal.

<sup>3</sup>In this table, *DN* abbreviates *Denver Nuggets* and *BC* abbreviates *Boston Celtics*.

One could reach a more efficient solution if the decision corresponding to frame T015 is delayed and the unifier commits to a value for P only after frame T100. The **wait** control specification is used to achieve this effect under certain conditions.

Wait specifies that the decision corresponding to a disjunction depends on the value of certain features. The syntax is as follows:

```
(alt TP (:wait P) ...)
```

This declaration indicates to the unifier that the alt TP depends on the value of the feature P. If when it is first met feature P is instantiated, then the alt is evaluated normally. If however P is not yet instantiated, the whole disjunction is delayed: it is put on hold, on an agenda. The rest of the grammar is evaluated, and periodically, the unifier checks the agenda to determine if one of the frozen alts can be awakened.

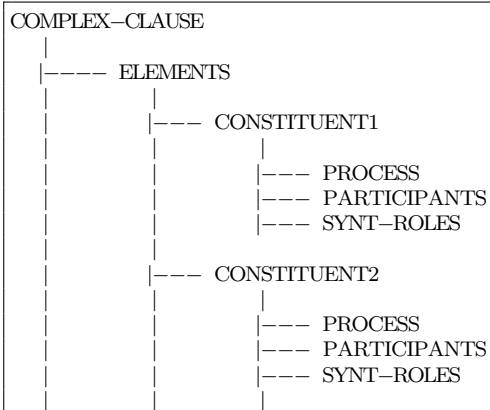
The following example illustrates a possible use of **wait**:

```
;; In conjunction: do verb ellipsis?
;; If the same verb is used in both conjuncts
;; you can ellide it in the second conjunct:
;; --> John ate the burger and Mary the fish.
;; The ellipsis is done by adding a GAP feature to the second verb.
;; Wait until verbs for both conjuncts have been selected.
;; This alt does NOT determine what the verbs should be.
(alt verb-ellipsis (:wait {constituent1 process lex}
                          {constituent2 process lex}))

(((cat clause)
  ({^ constituent1 process lex} {^ ^ ^ ^ constituent2 process lex})
  ({^ constituent2 process gap} yes)
  (verbal-ellipsis yes))
 ((verbal-ellipsis no))))
```

This fragment is extracted from grammar gr10. It is part of the branch of the grammar dealing with conjunction. This alt implements the decision whether to use a verb ellipsis in a conjunction of two clauses. The verb in the second conjunct can be elided if it the same lexical entry as the verb in the first conjunct. This match is checked by testing that the constituent1 process lex path can be unified with the constituent2 process lex path. When it can be unified, the constituent2 process is enriched with the feature (gap yes), which the linearizer will interpret by skipping the verb in the final sentence.

Now remember that the standard control flow in FUF is top-down breadth-first in the tree of constituents. A conjunction of clauses is made up of the following constituents:



When the branch of the grammar specifying the conjunction is traversed, the choice of the lexical item that will realize the process is still three levels down. The decision whether to use an ellipsis concerns only the conjunction grammar, so it must be located in the conjunction branch, not in the branch of the grammar deciding on lexical choice for the process. But it depends on this later decision - we do not want to influence the choice of the verbs in the conjunct by the fact that they are conjoined. So the decision is instead delayed until the verbs in both conjuncts have been determined.

Note that if the input already specifies the verbs, then the decision can be evaluated right away. It is only delayed when needed.

There can be complications when using `wait` when two decisions wait for each other, in a deadlock configuration. I have not yet encountered such situations in practice but in such cases, a combination of `wait` and `bk-class` could ensure that the deadlock can be broken in an efficient way.

### 12.4.1 Wait and Constituent Traversal

The delaying of disjunctions interacts in a complex manner with the way the constituent structure is traversed. Recall that the traversal proceeds as follows: the toplevel constituent is unified with the grammar, then at the end of the unification (before determination time), the `cset` of the constituent is identified, and every descendant specified in the `cset` is traversed, in order. The structure is then traversed breadth first.

As explained in Sect.5.7.1, there are two ways to specify a `cset` in a constituent in FUF: implicitly and explicitly. The `cset` is explicitly specified if there is a complete `cset` feature found in the constituent. Otherwise it is implicit. The implicit `cset` is found by using two heuristics: first it is assumed that any path mentioned in a pattern is also a constituent; second, any feature which contains a sub-feature of the form `(cat xx)` is also assumed to be a constituent.

To override these heuristics, the grammar writer can use *incremental cset specifications*. There are indeed two types of `cset` specifications: complete and incremental. A complete specification is of the form `(cset (c1 ... cn))` and exhaustively identifies all the subconstituents of the current constituent. An incremental specification is of the form: `(cset(+a1...am)(-d1...dp))`

An incremental `cset` specifies that all the `ais` should be added to the implicit `cset` and all the `djs` should be removed from it. The possibility of specifying the constituent structure incrementally is a great convenience practically since it allows the grammar writer to give partial information on constituency in different regions of the grammar.

Incremental `cset` specifications may affect delayed disjunctions. The following example illustrates this situation:

```
((alt (:wait {^ manner conveyed})
  (((manner ((conveyed any)))
    (manner ((conveyed adverb)
              (cat adverb) ...))
    (cset ((+ manner))))))))
```

In this case, the decision to express the manner constraint by using an adverb is delayed, leaving a chance for the verb to express it as a connotation. Thus, in the first pass through the grammar, this `alt` is not evaluated. At the end of the processing of the clause constituent, the `cset` is determined. The implicit `cset` is computed and all the incremental `cset` specifications are added on it. At this point, the `manner` feature is not identified as a constituent, since there is no `cat` specified for it and no incremental `cset` has added it explicitly. The constituent structure is then traversed. At the end of this traversal, the delayed `alts` on the agenda are processed. The `manner` `alt` is thus unblocked, and the second branch is eventually selected. An adverb constituent is then added at the clause level. When the `alt` is completely processed, this adverb constituent should be unified with the adverb grammar as any subconstituent would be, if it had been processed without being delayed. Unfortunately, the constituent structure was already computed before the delayed `alt` was unblocked. And the adverb was not part of the `cset` at this time. So the adverb is not added to the final constituent structure of the clause.

I have extended the determination process to re-check the constituent structure whenever an `alt` has been unblocked during determination, thus allowing new constituents to be added by delayed `alts` without restriction. The determination process performs the following task: first, a delayed `alt` is selected (the first one that was put on the agenda is chosen) and its evaluation is forced. Next, unification starts again to evaluate all the delayed `alts` which have been unblocked by the evaluation of the forced `alt`. At the end of this unification step, determination starts again. This loop continues until the agenda is empty. At this point, the constituent structure is recomputed and retraversed from the top of the total FD down, in breadth-first order. All constituents which have already been unified are skipped, all new constituents identified by the recomputation of the `csets` at the end of the agenda evaluation are unified again. After this unification step, determination starts again from scratch (since the evaluation of the new constituents may have added delayed `alts` on the agenda), until the agenda is empty and all constituents have been unified.

This top-down traversal in several passes integrates nicely the regular top-down control regime with the benefits of the dynamic re-ordering of constraints provided by `wait`. Note that this traversal in several passes does not translate in reduced efficiency: it is only triggered when needed (that is, when decisions had to be delayed), and

the traversal/recomputation of the constituent structure is a very efficient operation. As measured, efficiency is overall drastically improved with **wait**.

## 12.5 Conditional-Evaluation with Ignore

The **ignore** control annotation allows the grammar writer to evaluate certain decisions only under certain conditions. The idea is to just ignore certain decisions when there is not enough information, there is already enough information or there is not enough resources left. The 3 cases correspond to the annotations:

```
(alt (:ignore-when <path> ...) ...)
(alt (:ignore-unless <path> ...) ...)
(alt (:ignore-after <number>) ...)
```

**Ignore-when** is triggered when the paths listed in the annotation are already instantiated. It is used to check that the input already contains information and the grammar does not have to re-derive it.

**Ignore-unless** is triggered when a path is not instantiated. It is used when the input does not contain enough information at all, and the grammar can therefore just choose an arbitrary default. A real example of the use of **ignore-unless** is given below.

**Ignore-after** is triggered after a certain number of backtracking points have been consumed. It indicates that the decision encoded by the disjunction is a detail refinement that is not necessary to the completion of the unification, but would just add to its appropriateness or value. (IGNORE-AFTER IS NOT IMPLEMENTED IN FUF5.2).

One characteristics of these annotations is that their evaluation may depend on the order in which evaluation proceeds. Since this order is not known to the grammar writer, their use can be delicate. To prevent unpredictable variations, the **ignore** annotations should be used in conjunction with **wait**, since **wait** establishes constraints on when a decision is evaluated. Therefore a **wait** annotation has priority over an **ignore** when both occur in the same alt. Adding a **wait** annotation can often ensure that the **ignore** annotation will only be considered when it is relevant.

The following example is extracted from grammar gr10. It is actually the same case of verb ellipsis in conjunction presented in the section on **wait** above. In this decision, we also want to specify that the whole decision of whether to use ellipsis or not is only relevant in the case of a conjunction of clauses. This is expressed by addition the **ignore-unless** annotation as shown in the following figure:

```
;; In conjunction: do verb ellipsis?
;; This decision ONLY applies to CLAUSES.
;; Ignore it in conjunctions of NPs and other cats.
(alt verb-ellipsis (:wait process)
  (:ignore-unless ((cat clause)))
  (((cat clause)
    ({^ constituent1 process lex} {^ ^ ^ constituent2 process lex})
    ({^ constituent2 process gap} yes)
    (verbal-ellipsis yes))
   ((verbal-ellipsis no))))
```



## Chapter 13

# Tracing and Debugging

### 13.1 What it Means to Debug a FUF Program

When using FUF, a grammar developer is programming in the FUF programming language. The grammar is a program. The input FD is the input to the program. FUF is the interpreter and sentences are the output of the execution of a grammar on inputs. In this framework, when something “goes wrong,” the grammar developer must debug his grammar. The main sources of bugs found when developing FUF programs are:

- The input is not well formed (it is not a valid FD).
- The grammar is not well formed (it is not a valid FD).
- The input does not have the structure expected by the grammar (too flat or too deeply nested).
- The constituent structure is badly specified in the grammar (causing either too many constituents to be generated, or not enough, or infinite recursion in the constituent traversal).
- The ordering patterns are not correctly specified or two pattern specifications do not unify.
- The appropriate morphological features are not passed to the leaf-constituents, preventing the morphology module from performing the required inflections.
- Relative paths are not pointing to the place they were intended to. An ambiguous relative path is not correctly resolved.
- Given/any fail when they should not: given can interact poorly with wait, and should be replaced by any, any can cause heavy backtracking and should be replaced by given.
- Types are not defined as expected, or there is an unwanted interaction between types.
- FSET declarations are too rigorous.
- Wait and/or bk-class are not used correctly.
- An indexed feature is not found in an indexed alt.

The main problem when debugging a FUF program is to identify what caused a failure to occur. We introduce the following terminology to discuss the various debugging tools available in FUF: a failure occurs whenever the unifier attempts to unify a simple leaf symbol with a different, incompatible leaf. Failures trigger backtracking. An unexpected failure is a failure that the programmer did not expect. The initial failure is the first unexpected failure to occur during an unification. The distinction between regular failure and unexpected failure is of course subjective, but it is useful, because there are usually many failures that occur even if there are no bugs in the grammar. This happens whenever a non-indexed alt is traversed. Branches are tried in order until the appropriate branch is found. When failure occurs during the test of the initial branches, this is not the trace that “something is going wrong”. So the main problem of the FUF debugger is to, as quickly as possible, identify the initial failure - and to filter out all the irrelevant expected failures.

This task is made difficult because FUF backtracks a lot, and if the initial failure is missed, a lot of subsequent “unexpected failures” will occur when wrong branches are tried upon backtracking (in a sense, everything that happens after the initial failure should be disregarded, as the unifier is engaging on a wrong course).

This chapter presents the tracing tools available in FUF and provides advice on how to use them to quickly identify the initial failure.

## 13.2 Checking the Validity of FDs and Grammars

Before tracing the unification, it is important to check that both the input and the grammar are valid FDs. The following functions perform this checking:

```
(FD-SYNTAX fd)
;; check that a Lisp expression FD is a well-formed FD.

(FD-P fd)
;; check that a well-formed FD does not contain inconsistencies
;; (that is, (u fd nil) is not :fail, or in other terms, the FD
;; does not contain contradictory features, such as ((a 1) (a 2))).

(GRAMMAR-P grammar)
;; check that a Lisp expression grammar is a well-formed grammar.
```

When you suspect that your current input/grammar do not work, check first that they are syntactically valid using these functions.

## 13.3 Fine Tuning Tracing: Overview of FUF Tracing Functions

Several dimensions characterize the activity of the FUF unifier:

- Where in the *input* is the unification proceeding.
- Where in the *grammar* is the unification proceeding.
- How *important* is the current action of the unifier.
- What *stage* of the unification is proceeding.

It is possible to tailor the tracing behavior of FUF according to each of these dimensions. In general, FUF can output a tracing message whenever it takes an action, such as adding a feature to the total FD, selecting a branch in the grammar, backtracking because of a failure, expanding a cset and moving in the constituent tree traversal, freezing and unfreezing goals etc. Outputting all the possible trace messages is always overwhelming, and provides too much text to be useful. So the challenge of the FUF debugger is to fine-tune the tracing system to produce just enough information to locate the bugs in the grammar and/or in the input FD.

The grammar developer indicates what portions of the grammar must be traced: the grammar is traced, not the unifier. Therefore, to trigger tracing, one must put directives into the grammar. At the Lisp level, and for a given grammar including tracing directives, traces can be switched on or off by the following functions:

```
;; I GENERAL TRACING CONTROL

(trace-on) ;; enable all trace messages to be output.
(trace-off) ;; disable all trace messages to be output
(trace-bp &optional (frequency 10))
;; Output a dot for every frequency backtracking points.
;; Useful for long computations to get a feeling of what's happening.
;; Works even if (trace-off).
%break% ;; allows the insertion of break points in the grammar.
(trace-level level)
;; Determines detail level of trace to be printed. The following
;; levels are defined:
;;          00: feature level action
;;          05: unimportant alt-level action
;;          10: alt-level action - branch number trying
;;          12: demo messages
;;          15: freeze, ignore, bk-class
;;          20: constituent level action
;;          30: important failure - end of alt

;; II TRACING FLAGS MANAGEMENT: ENABLE & DISABLE

(all-tracing-flags &optional (grammar *u-grammar*))
```

```

;; return the list of all tracing flags defined in grammar.
(trace-disable flag)
;; disable flag.
;; Everything works as if flag was not defined in the grammar.
(trace-enable flag) ;; re-enable a disabled flag.
(trace-disable-all) ;; disable all flags.
(trace-enable-all) ;; re-enable all flags.
(trace-enable-alt alt-name :expansion t :grammar g)
(trace-disable-alt alt-name ...)
;; Enable all tracing flags defined under a def-alt or
;; def-conj. If expansion is nil, does not expand the
;; sub-def-alt.
(trace-disable-match string)
;; disable all flags whose names contain string.
(trace-enable-match string)
;; re-enable all flags whose names contain string.

;; III CONTROL OF SPECIFIC ACTIVITY TRACING

(trace-determine :on t|nil) ;; enable tracing of determine stage or not.
(trace-category :all|cat|(cat1...catn) t|nil)
;; enable tracing of categories or not.
(trace-bk-class t|nil) ;; list special messages concerning bk-class
(trace-wait :on t|nil) ;; list special messages concerning goal freezing
(trace-cset :on t|nil) ;; trace cset expansion
(trace-alts :on t|nil) ;; detailed tracing of all alts, even unnamed.
(hyper-trace-category cat :status t|nil)
;; trace category with printing of full constituent
;; before unification of constituents of the traced cats.

```

## 13.4 Identifying Possible Bugs: Trace-bp

To identify whether FUF enters into unbounded backtracking (which is the sign of a bug in the grammar or in the input), it is useful to first monitor how many backtracking points are being used. This can be achieved by using the **trace-bp** (trace backtracking points) function. This is the only tracing function which can be activated even if the general tracing system is turned off (with the function **trace-off**). The effect of the function is to output a dot on the screen every *n* backtracking points, where *n* is 10 by default and can be changed by giving an argument to **trace-bp**. The form of the function is:

```

(TRACE-BP &optional n)

(TRACE-BP nil): turns off printing of .

```

If the string of dots becomes longer than expected, then it is worth stopping the unifier under the suspicion that something is going wrong, and to start the tracing system with **trace-on**.

## 13.5 Levels of Tracing

In the attempt to reduce the amount of trace messages and finding the original source of failure, the most useful tool is the **trace-level** function, which filters tracing messages according to their importance. The following levels of importance are predefined:

- 30: important failure - end of alt
- 20: constituent level action
- 15: freeze, ignore, bk-class
- 12: demo messages
- 10: alt-level action - branch number trying
- 05: unimportant alt-level action
- 00: feature level action

The function `trace-level` sets the minimum level of tracing messages that can be printed. Thus, calling `(trace-level 20)` insures that only important messages of level 20 and 30 will be printed.

The levels are defined as follows:

- 30: An “important” failure in unification occurs - that is, all the branches of an alt have been tried and failed. This in general means that the input is not compatible with the grammar. It does not force immediate failure of the overall unification process, because some backtracking options may still be open, but it is a strong indication that something wrong is going on. In general, there is a high probability that the *first* level 30 failure message is the initial failure (cf. p.63).
- 20: Constituent Level Action - traces the constituent tree traversal of FUF. Whenever the grammar is re-accessed to unify a new constituent, a tracing message is printed. This is useful to follow the progression of the unifier through the total FD. In general, the traversal is a top-down breadth-first expansion of the constituent tree.
- 15: Control messages: The dominant control strategy of FUF is the top-down constituent traversal. Fine-tuning of this strategy is, however, possible, using the `bk-class`, `wait` and `ignore` directives. These constructs are described in more detail in Section 12.
- 12: Demo messages: The alt construct allows the grammar writer to output “demo messages” when an alt is entered. These messages are considered of level 12.
- 10: Alt-level action - branch number trying: when an alt is tried, branches are tried successively, in order, randomly (for a ralt) or directly (for an indexed alt). A tracing message is output each time a branch is tried, indicating the name of the alt (therefore the position of the unifier within the grammar) and the number of the branch within the alt.
- 05: unimportant alt-level action: when an alt is indexed, certain tracing messages are printed to document the index search.
- 00: feature level action: each time a feature is added to the total FD by the unifier, a tracing message is printed. This is the absolute lowest level of tracing and results in unending seas of messages.

In general, when debugging, start by setting `(trace-level 30)`, this provides most of the time directly the location of the initial failure.

## 13.6 Tracing of Alternatives and Options

To follow the unifier as it proceeds through the grammar, the most useful trace is generated by giving a name to an alternative of the grammar. It is done by adding an atomic name after the keywords `alt`, `ralt` or `opt` in the grammar:

```
((alt PASSIVE
  (
    ;; branch 1 of alt passive
    ((verb ((voice passive)))
      (prot none))

    ;; branch 2 of alt passive
    ((verb ((voice passive)))
      (prot any)
      (prot ((cat np)))
      (by-obj ((cat pp) (prep ((lex "by"))) (np (^ prot))))
      (pattern (dots verb by-obj dots))))

    ;; body of alt passive (common to all branches)
    (verb ((cat verb-group)))
    ...))
```

Here, this fraction of the grammar has been marked by the directive: `(alt PASSIVE ...)`. (An equivalent notation is `(alt (:trace PASSIVE) ...)`.) The effect will be that all unification done subsequently will be traced, producing the following output:

```
--> Entering ALT PASSIVE
--> Trying Branch #1 in ALT PASSIVE:
--> Fail on trying (prot none) with
      (prot ((nnp ((n ((lex boy)))))))
--> Trying Branch #2 in ALT PASSIVE:
```

...

If a traced alternative is found later in the grammar, the level of indentation will increase. If the level of indentation decreases, that means a whole (**alt** ...) has failed. It is indicated by the output:

—> Fail on ALT PROT.

The possible messages printed when the grammar is traced are:

```
;; Move in the alternatives:
;;   ENTERING ALT f: BRANCH #i
;;   FAIL IN ALT f
;;   When the alt is indexed:
;;   ENTERING ALT f — JUMP INDEXED TO BRANCH #i INDEX-NAME
;;   NO VALUE GIVEN IN INPUT FOR INDEX INDEX-NAME — NO JUMP
;; For options:
;;   TRYING WITH OPTION o
;;   TRYING WITHOUT OPTION o
;; Regular unification:
;;   ENRICHING INPUT WITH s AT LEVEL l
;;   FAIL IN TRYING s with s AT LEVEL l
;; Pattern unification:
;;   UNIFYING PATTERN p with p
;;   TRYING PATTERN p
;;   ADDING CONSTRAINTS c
;;   FAIL ON PATTERN p
;; Unification between pointers to constituents:
;;   UPDATING s WITH VALUE s AT LEVEL l
;;   s BECOMES A POINTER TO s AT LEVEL l
;;   UPDATING BOTH PATHS TO A BOUND
```

## 13.7 Local tracing with boundaries

If you want a more focused tracing, you can put anywhere in the grammar a pair of atomic flags whose first character must be a "variable **\*trace-marker\***"). All the unification done between the 2 flags will be traced, and will produce the same messages as usual.

```
;; branch 2 of alt passive
((verb ((voice passive)))
 (prot any)
 %by-obj%
 (prot ((cat np)))
 (by-obj ((cat pp) (prep ((lex "by")) (np (^ prot))))
 %by-obj%
 (pattern (dots verb by-obj dots)))
...
```

All the unification done between the 2 flags You furthermore will have a message:

```
Switching local trace flags on and off:
    TRACING FLAG f
    UNTRACING FLAG f
```

You generally want to have only small portions of the grammar put between tracing flags.

### 13.7.1 Special Flags **%trace-on%** and **%trace-off%**

The special tracing flags **%trace-on%** and **%trace-off%** are predefined and have the effect of temporarily turning all tracing on and off. They are used as all other local tracing flags, as shown in the following example:

```
((alt PASSIVE
 (
  ;; branch 1 of alt passive
  ((verb ((voice passive)))
   (prot none))
```

```

;; branch 2 of alt passive
((verb ((voice passive)))

%trace-off%                ;; Turn off tracing of details in this branch

(prot any)
(prot ((cat np)))
(by-obj ((cat pp) (prep ((lex "by")))) (np (^ prot)))
(pattern (dots verb by-obj dots))

%trace-on%                ;; Turn tracing back on

)))

;; body of alt passive (common to all branches)
(verb ((cat verb-group)))
...)

```

**%Trace-off%** is generally used to remove tracing information from a region of the grammar that has already been debugged but still leaving the alt traversal tracing messages on. The same effect is most of the time achieved by using the **trace-level** function, but **%trace-off%** allows finer tuning of the tracing behavior of FUF.

### 13.7.2 The Special Tracing Flag **%break%**

The special tracing flag **%break%** is used to trigger a break in the execution of FUF. When found by FUF, a Lisp continuable break is triggered. It is then possible to examine the current state of the unification using the Lisp debugger. Using the Lisp debugger, it is then possible to either resume unification or abort it. Within the debugger, the total fd can be inspected (and modified) by using the function **path-value** and **set-path-value**. A typical session is shown below:

```

(setq *u-grammar*
  '((alt (
    ((cat clause) (pattern (s v o))
      (s ((cat np)))
      (v ((cat v)))
      (o ((cat np))))
    ((cat np) %break%)
    ((cat v) %break%))))))

;; The np and v branch are not yet written - a break is inserted
;; The developer can then insert new values manually during unification.

LISP> (uni '((cat clause)
            (s ((lex "John"))))
        (v ((lex "like"))))
        (o ((lex "Mary"))))

>=====
>STARTING CAT CLAUSE AT LEVEL {}
>=====
>Expanding constituent {} into cset ({O} {V} {S}).
>
>=====
>STARTING CAT NP AT LEVEL {O}
>=====

Break: Break in grammar

Restart actions (select using :continue):
0: return from break.
{1c} FUG5 23> (path-value {o})

((LEX "Mary") (CAT NP))
{1c} FUG5 24> (path-value {n})

((LEX "John") (CAT NP))

{1c} FUG5 25> :cont

```

```

>Constituent {O} is a leaf.
>
>
>=====
>STARTING CAT V AT LEVEL {V}
>=====
Break: Break in grammar

Restart actions (select using :continue):
  0: return from break.
{1c} FUG5 26> :reset

FUG5 27>

```

The behavior of the debugger depends on which version of Common Lisp you are using. The example shown here is run under Franz Inc's Allegro Common Lisp. Consult your Lisp manual to find out the details of how to resume processing (the `:cont` command in ACL) and abort processing (the `:res` command in ACL). Another source of variation is whether the `{}` notation is recognized within the debugger or not. This notation is implemented using macro characters in Lisp. Macro-characters are recognized in ACL's debugger but not in Lucid Common Lisp's implementation. For that reason, the function `path-value` accepts as parameter either a path or simply a list of attributes.

The function `path-value` returns the value of a path in the current total FD. It is useful to inspect the current value of the total FD. The function `set-path-value` is also defined to change a value within the total FD. Note that its use is highly dangerous.

```

(path-value path-or-list)  Return the value of path in the current total FD.

(set-path-value path-or-list FD)  Set the value of a path in the current
                                total FD to FD.

;; Examples:

(path-value {process v})
(path-value '(process v))  ;; equivalent form when the {} notation is not
                           ;; recognized

(set-path-value {process v} '((lex "take")))

(set-path-value {process v} (u (path-value {process v}) '((tense past))))
                           ;; u performs a simple unification between fds.

```

For more general access, the total FD is accessible in the special variable `*input*` and it can be modified in any possible ways. But if you do follow this way, there is a high probability that the unification will not be able to proceed normally. Note that there is no way to easily remove a conflation from the total FD using only `path-value` and `set-path-value` because `path-value` always follows the paths until a non-path value can be returned. The following example illustrates this limitation:

```

(setq *input* '((a {b})
               (b ((b1 1)))))

(path-value {a}) --> ((b1 1))

(set-path-value {a} '((b1 2)))

(path-value {}) --> ((a {b}) (b ((b1 2))))

;; Cannot just with path-value and set-path-value remove the conflation
;; between a and b (a {b}).

```

One last word of caution: when using `set-path-value`, relative paths are made absolute before being inserted relative to the path of insertion given as parameter. Since the relative path notation can be ambiguous (cf p.17), this can, under circumstances where there is an ambiguity, create unexpected results.

## 13.8 The trace-enable and trace-disable Family of Functions

In general, a grammar is defined in a file, that you load in your Lisp environment. The tracing flags are defined in that file after the alts and opts or as local flags. When you develop a grammar, you want to focus on different

parts of the grammar. In order to do that, you can selectively enable or disable some of the flags defined in the grammar.

The function `all-tracing-flags` returns a list of all the flags defined in the grammar. You can then choose to enable or disable all the flags, only a given flag, or all flags whose name matches a given string.

When a flag is disabled, everything happens as if the flag was not defined at all in the grammar. Note that you cannot create a new flag in the grammar by using these functions. You can simply turn on and off existing flags. It is therefore a good idea to define all the possible flags in a grammar and to adjust the list of enabled flags from within lisp.

When you use the `def-alt` and `def-conj` notation (cf. Section 6), the functions `trace-enable-alt` and `trace-disable-alt` can be used to enable and disable all the tracing flags appearing under a given `def-alt` or `def-conj` construct. The `expansion` keyword determines if the enabling/disabling only concerns the tracing flags appearing directly in the `def-alt` construct, or also all the flags appearing in the expansion of the construct.

## 13.9 The :demo directive

Reading traces from the unifier is a particularly tedious task. The main problem is that the messages generated by the program are very similar to each other. The `:demo` directive allows the grammar writer to bring some variety to these messages. A demo-message can be used to output a specific message during the trace of the program when entering an `alt` (or a `ralt`) construct. The syntax is indicated by the following figure:

```
(alt voice (:demo "Is the voice active or passive?")
  (((voice active)
    ...)
   ((voice passive)
    ...)))
```

The message will be printed in the trace of the program (only if the tracing flag `voice` is enabled) as shown:

```
--> Entering alt VOICE
      Is the voice active or passive?
--> Trying branch #1
      ...
```

Note that the message is indented in the stream of trace messages. Such messages allow the grammar writer to put some semantic information into the trace messages, so that the whole stream of messages can be more easily interpreted.

In addition to its position at the top of an `alt` construct, a demo-message can be embedded anywhere in a grammar by using the `control-demo` function. `Control-demo` must be used within a `control` pair and produces an indented demo message in the trace stream. The following example illustrates its use: `control-demo`

```
(alt from-loc (:demo "Is there a from-loc role?")
  (((from-loc none)
    %TRACE-OFF%
    (control (control-demo "No from-loc"))
    %TRACE-ON%)
   ((from-loc given)
    %TRACE-OFF%
    (control (control-demo "From-loc is here"))
    %TRACE-ON%)))
```

*;; Tracing output:*

```
--> Entering alt FROM-LOC
      Is there a from-loc role?
--> Fail with branch #1
--> Entering branch #2
      From-loc is here
--> Success with branch #2
```

NOTE: The `control` pair containing a `control-demo` call should be put within a `%TRACE-OFF%` - `%TRACE-ON%` pair to avoid the printing of system trace messages regarding the control pair.



NOTE: The demo message string is passed to the `format` common-lisp function, and can therefore contain formatting characters accepted by that function (*e.g.*, `~` or `~%`). Refer to your CommonLisp manual for details.

## 13.10 Tracing of Specific Stages of the Unification

The following group of function enable or disable tracing of specific stages of the unification process:

```
(trace-determine :on t|nil)
(trace-category :all|cat|(cat1...catn) t|nil)
(trace-bk-class t|nil)
(trace-wait :on t|nil)
(trace-cset :on t|nil)
(trace-alts :on t|nil)
(hyper-trace-category cat :status t|nil)
```

### 13.10.1 Trace-determine

The determination stage checks that no ANY constraint is left unsatisfied at the end of the unification stage, and that all TEST constraints are satisfied. In addition, it checks that no further CSET traversal is required to restore after frozen constraints have been thawed or forced.

When `trace-determine` is on, these checks generate tracing messages. The specific messages are:

```
FAIL: Found an ANY at level <path>
Current Sentence: ....

TEST succeeds: <expr> at level <path>

Fail in testing <expr> at level <path>
```

The current sentence is only printed when FUF is called from the toplevel function `uni` whose function is to print a sentence. In other uses of FUF, the current function is not generated. This message is of level 30 (highest level).

In addition, each time the determination stage is reached, a line of statistics on backtracking is printed, of the form:

```
{Used 119 backtracking points - 26 wrong branches - 15 undos}
```

Three pieces of information are provided:

1. The total number of backtracking points used so far. This can be interpreted as the number of “questions” FUF asked to the grammar, or how many times a branch had to be chosen in an alt construct.
2. The number of “wrong branches” - that is, how many times FUF picked up a branch and started unification to finally undo this branch upon backtracking. This can be interpreted as the number of “wrong guesses” made by FUF.
3. The number of “undos” - that is, how many features were added to the total FD while traversing “wrong guesses” and were later removed from the total FD upon backtracking. This gives an indication of how “deep” FUF went into wrong branches before realizing that they led to failure.

When, because of the interaction between wait and the constituent traversal (cf Section 12.4.1) a new cycle of unification must be started after determination, a new line of statistics will be printed. These lines of statistics are NOT printed if the keyword parameter `non-interactive` is set to true when calling the top-level functions of FUF (for example, as in `(uni fd1 :non-interactive t)`).

These statistic lines are printed indicate whenever a unification cycle ends and a determination cycle starts. They therefore provide important information on how unification is proceeding.

### 13.10.2 Trace-Category and Hyper-trace-category

Trace-category is useful to follow the traversal of the constituent structure as it is performed by FUF. When a category C is traced, the following message is printed whenever a constituent of category C is unified with the grammar:

```
>STARTING CAT ADJ AT LEVEL {SYNT-ROLES SUBJ-COMP HEAD}
```

The function `hyper-trace-category` provides more detail: if category C is “hyper-traced”, the same message is printed, and in addition the whole constituent is printed:

```
>STARTING CAT ADJ AT LEVEL {SYNT-ROLES SUBJ-COMP HEAD}
```

```
>CONSTITUENT {SYNT-ROLES SUBJ-COMP HEAD} =  
((CAT ADJ) (CONCEPT {SYNT-ROLES SUBJ-COMP CONCEPT})  
 (POLARITY {SYNT-ROLES SUBJ-COMP POLARITY})  
 (LEX {SYNT-ROLES SUBJ-COMP LEX})))
```

The functions are used as follows:

```
(TRACE-CATEGORY c | :all | (c1 ... cn) &optional t | nil)  
(HYPER-TRACE-CATEGORY c | :all | (c1 ... cn) &optional t | nil)  
  Trace or hyper-trace a given category, or all categories  
  (if parameter is :all) or a list of categories.  
  If a second parameter nil is added, the category or  
  categories are untraced.
```

### 13.10.3 Trace-Cset

The `trace-cset` function is used to monitor constituent traversal. Each time FUF finishes unifying a constituent, it computes the cset of the result of the unification, applying the rules described in Section 5.7. The constituents are then traversed breadth-first. The function `trace-cset` triggers messages of the following form at the end of the unification of each constituent:

```
>Expanding constituent {} into cset ({SYNT-ROLES SUBJ-COMP}  
                                     {PROCESS}  
                                     {SYNT-ROLES SUBJECT}).
```

```
>STARTING CAT AP AT LEVEL {SYNT-ROLES SUBJ-COMP}
```

If the constituent has an empty cset, it means it is a leaf in the constituent structure. In that case, the following message is printed:

```
>Constituent {SYNT-ROLES SUBJ-COMP HEAD} is a leaf.  
>Constituent {PROCESS EVENT} is a leaf.  
>
```

### 13.10.4 Trace-BK-Class

The interpretation of the BK-class annotation has been introduced in Section 12.3. BK-class is used to allow intelligent dependency-directed backtracking. Each choice point can be marked to belong to a certain backtracking class (bk-class), and certain classes of paths can be declared to belong to corresponding bk-classes. Upon failure, the address of failure is checked. If it does not belong to a declared bk-class, regular chronological backtracking is started. If it does belong to a bk-class, then backtracking continues up to the latest choice points that belong to the same bk-class.

When tracing this behavior, the following conditions are monitored: (1) when a special failure address is met and how high up the intelligent backtracking progresses and (2) what is the latest address of failure. The following messages are printed in each case:

```
>BK: Special path <path> caught by alt <alt> of class <c> after <n> frames.
>BK: Switch from <path1> to <path2>
```

The first message is emitted whenever intelligent backtracking occurs. The second one is emitted whenever the current address of failure is modified following the rules discussed on page 58.

### 13.10.5 Trace-Wait

The wait annotation is used to allow goal delaying during the evaluation of a grammar. The idea is to wait until enough information is available before trying to choose between the branches of an alt. The grammar writer indicates which information is requested before the evaluation of an alt can start using the wait annotation and a sequence of path expressions which point to the features which must be instantiated with the requested information.

When enough information is already present, the alt is evaluated as usual. When there is missing information (some of the features are not yet instantiated), the alt is frozen (delayed). Frozen alts are stored on an agenda, which keeps track of the decisions which remain to be made in the future. At regular intervals (in FUF, whenever a new choice point is met), the agenda is checked, and if enough information has been gathered since the time a decision has been frozen, the decision is thawed, and evaluated immediately. After the evaluation of the thawed alt, control proceeds where it was interrupted.

In addition, at the end of the unification stage, the determination stage checks if any decision is still on the agenda. Such a situation is reached if not enough information could be gathered anyway to allow the evaluation of the frozen alts. In this case, the frozen alts of the agenda are “forced” and evaluated, even though some of the requested information is missing.

Finally, there is one more configuration of control decisions which concerns the goal delaying behavior of FUF: as explained on page 62, wait annotations have priority over ignore annotations to insure that ignore annotations are only considered when enough information is available for them to make sense. When thawing a frozen alt from the agenda, after enough information has been gathered, the ignore annotations are immediately checked. If they do match, then the whole frozen alt is immediately ignored.

Corresponding to each of these configurations, the trace system emits the following messages. In every case, a unique agenda identifier (integer number) is assigned to each frozen alt:

```
>Freezing alt <alt>: waiting for <path> {agenda <n>}
>Thawing {agenda <n>}: Restarting at level <path>
>Forcing {agenda <n>}: Restarting at level <path>
>Ignoring {agenda <n>}
```

### 13.10.6 Trace-Alts

The function `trace-alts` systematically monitors traversing of alts in the grammar, even if the alts are not traced, and outputs a message whenever a new branch is tried. It is best used in conjunction with a `trace-disable-all` setting, to follow uniquely alt traversal. The form of the output is as follows:

```
LISP> (trace-disable-all)
LISP> (trace-alts)
LISP> (uni t1)

=====
>STARTING CAT CLAUSE AT LEVEL {}
=====

>Fail in alt VOICE
>Fail in alt VOICE-NORMAL
```

```

>Fail in alt VOICE-NORMAL
>Fail in alt VOICE-NORMAL
>Fail in alt DATIVE-MOVE-DEFAULT
>Fail in alt DATIVE-MOVE-DEFAULT
>Fail in alt SUBJECT-SUBCAT
>Fail in alt SUBJ-COMP-CAT
>Fail in alt :ANONYMOUS
>Fail in alt :ANONYMOUS
>Expanding constituent {} into cset ({SYNT-ROLES SUBJ-COMP}
                                     {PROCESS}
                                     {SYNT-ROLES SUBJECT}).
>

```

Each time a branch in an alt fails, the message “fail in alt X” is printed. If the alt is not traced, the name `:anonymous` is printed instead. This is useful to find possible errors even in places which are not traced in the grammar. In general, `trace-alt` should only be used in last recourse.

## 13.11 Some Advice on FUF Debugging

This section provides some pragmatic advice on how to use the tracing system of FUF based on the experience gathered while developing large grammars. It lists common sources of confusion, warns against the misfeatures of the tracing system, common bugs and some successful bug tracking tricks.

### 13.11.1 Syntax Errors

The first source of bugs, often incomprehensible ones, is syntax errors, either in the input FD or in the grammar. So the first precaution is to check both inputs and grammars regularly. Use functions `FD-P` and `GRAMMAR-P` for that purpose. In grammars, check especially the number of parentheses occuring after ALTs.

### 13.11.2 Semantic Errors

Semantic errors in the grammar can be trivial typos or the sign of a bad design of the FD structure. In any case check the following points:

- Misspellings: wrong feature-name, wrong leaf-value-name
- Paths: in both forms, *i.e.*, `((f1 ((f2 ((f3 and {f1 f2 f3}`. Make sure that paths actually point to the location in the total FD that you expect.
- The most common source of error is to include the wrong number of up-arrows `^` in a path expression. Always double-check your relative paths, especially those appearing embedded in ALTs, `CONTROL`, `EXTERNAL` when the up-arrow notation can be quite counter-intuitive.
- Observe the structure of the unified graph before unification. To this end, use the function `uni-fd` and pretty-print its output as in:

```
(pprint (clean-fd (uni-fd input :grammar gr)))
```

This is often instructive. Use the function `top-gdp` to explore this output FD in detail.

- Draw a graph of the total FD with all features instantiated to help you check all your paths and levels of embedding.
- If you do not understand the current structure of the FD, insert a `%break%` in your grammar at some critical point, and use the function `path-value` to inspect the total FD during the time of the unification.
- It is risky to use absolute paths in general (SURGE does not include a single absolute path for example). Using an absolute path is a sign of desperation.

### 13.11.3 Expression of Negative Constraints

Negative constraints are used to limit the scope of acceptable FDs by the grammar and to force failure when certain FD configurations are met. The main tools in FUF for the expression of negative constraints are FSET and NONE. The main sources of confusion here are:

- FSET is too restrictive: you didn't plan on adding a feature, a new feature is not compatible with FSET.
- FSET must include explicitly all the features, including CAT and CSET.
- NONE is used in the wrong place.

### 13.11.4 Control

The overall flow of control of FUF is quite complex, and the trace messages do not make it very easy to follow. Before going on and suffering through the zillions of lines produced by the trace system, try to understand analytically what could go wrong - recheck the syntax of your fd and grammar, and recheck the structure of your FD and the structure your grammar builds. Only when this fails should you try to read the trace messages to understand what went wrong on a particular input. You should proceed in the following order:

1. Identify that there is a bug: do (trace-off) and (trace-bp). This will emit a dot every 10 backtracking points (10 is the default) and indicate how much effort FUF has invested on your input. For example, when dealing with SURGE, a clause that takes more than 300 backtracking points (30 dots on the screen) is either very complicated or contains a bug. When you think there is a bug, interrupt the unifier (generally use control-C).
2. Monitor the speed of appearance of the dots. In general, FUF operates in two modes: "forward" and "backward". Forward mode is when the grammar is traversed as expected, and every new feature falls into place. Backward mode is after the occurrence of the first unexpected failure. The unifier starts backtracking and tries every other branch in an unexpected manner. In general, these tried branches fail very quickly. So in backward mode, backtracking dots tend to be emitted much more quickly than in forward mode. When the dots start piling up, it's a good sign that FUF has entered backward mode, and that a bug has been found.
3. Once you suspect there is a bug, identify the first unexpected failure. The first attempt to do that is to set (trace-on) and (trace-level 30). This will print the most obvious candidates. If this does not work, basically, you're in trouble, and finding the bug will take time (sorry! I'm in the middle of the implementation of a graphical debugger for FUF which should help you beyond this stage).
4. The next step is to gradually lower the trace-level until you get a good understanding of where in the grammar is the source of the bug - moving to values 20, 15, 12, 10, 5 and 0 (trying trace-level 0 on a full grammar without disabling most of the tracing flags is a sign of desperation).
5. Once the approximate location of the problem is identified in the grammar, enable only the relevant tracing flags (those defined in the grammar around the problematic spot). This is achieved by using the functions (trace-disable-all) (trace-enable-alt name-of-suspected-alt).
6. From then on, try to understand what's happening. I find it convenient to run FUF within an Emacs buffer and to use the editor to search through the trace messages printed by FUF.
7. If you enable all tracing flags and set trace-level to 0 and you still cannot find a tracing message identifying a failure (of the form "FAIL in ...") this can be due to 2 problems:
  - (a) Either the failure is occurring in a region of the grammar which does not contain tracing flags. In this case, set (trace-alts) and check the messages "fail in alt :anonymous".
  - (b) Or the failure is due to a missing cat in the grammar. Often if you unify ((cat xxx)) with a grammar that has no branch for xxx, there is no failure message produced. Check your cats often.

The main points you should be looking for when debugging are:

- Endless backtracking: this can be found by setting trace-level 00 and looking for repetition in the flow of messages.
- Failure on ANY (especially in conjunction with recursion).

Failure on given (especially in conjunction with wait): remember that an alt that is annotated with (wait path1) can be forced and therefore evaluated even if path1 is not yet instantiated. Therefore, using (path1 given) in such a situation is dangerous. In this case, switch to any.

INDEX (especially double indexes and partial indexes). Do NOT factor together branches in an indexed alt. For example, the following index will not work (unification with ((a 1)) will fail):

```
(alt (:index a)
  (((a ((alt (1 2))))))
  ((a 3))))
```

WAIT (especially in conjunction with ignore and bk-class)

BK-CLASS: especially, since the bk-class declarations are persistent, do not forget to evaluate (clear-bk-class) when loading a new grammar. In case of doubt, evaluate (clear-bk-class) and re-evaluate all your (define-bk-class). In general, it is good to add the following prelude in all the files containing a grammar definition:

```
(eval-when (load eval)
  (clear-bk-class)
  (reset-typed-features))
```

IGNORE: check that you do not ignore alts too liberally.

Constituent Structure Definition:

- Wrong cset (especially interaction between cat and explicit cset). In case of doubts, use an explicit cset in your grammar. Avoid over-restrictive csets (with the = notation), they often cause failure of unification with further cset features in the grammar.
- Wrong pattern and interaction between pattern and cset when using implicit cset expansion.
- Csets can trigger infinite recursion.

Linearizer: The linearizer produces strings of the form “unknown cat: XXX” when a category unknown to the morphology module is found as a leaf in the constituent tree. This is often due to the following effect: all top-level unification functions have a :limit keyword option used to limit the number of backtracking point used on a single input. The default value for limit is 10,000. When unification stops after reaching this limit, the FD is sent to the linearizer, even though it has not been completely unified, and all constituents have not been completely expanded. In that case, you can get strings such as “unknown cat: NP” because an NP has not been expanded. To correct this problem, increase the limit.

Type Hierarchy Problems:

- Since the type declarations are persistent, do not forget to evaluate (reset-typed-features) when loading a new grammar. In case of doubt, use **draw-types** to verify the current state of type definitions. In general, it is good to add the following prelude in all the files containing a grammar definition:

```
(eval-when (load eval)
  (clear-bk-class)
  (reset-typed-features))
```

- Wrong **under** specification.
  - Wrong **define-feature-type**
  - Lack of **under** (when one wants only to test for the presence of a feature and not to enrich the total FD).
- Procedural types, test and control:
    - Wrong CONTROL/TEST function
    - TEST for CONTROL and vice-versa.
    - Wrong use of the path construct with a TEST/CONTROL (with the # notation).
    - Wrong EXTERNAL function.
    - Wrong procedural-type definitions.

## Chapter 14

# Manipulation of FDs as Data-structures

FDs can be viewed as a convenient data-structure. FUF provides a library of Lisp functions to manipulate FDs: accessors to extract a sub-fd from a total fd, insertion of an FD within a total FD, etc. These functions' main complexity lies in the interpretation of paths. For example, if you want to extract a sub-fd from a total fd, and this sub-fd contains path pointing outside the sub-fd up into the total fd, the extraction function should also extract these paths. In addition, relative paths in the sub-fd must be adjusted to maintain consistency.

Another class of functions helps in manipulating lists of FDs and FDs as lists (providing programmatically the same facilities as the  $\sim$ ,  $\hat{\sim}$  and  $\sim n$  notations discussed in Chapter 8 p.33).

In this chapter, the notion of *total FD* plays a critical role. A total FD is a self-contained FD, which appears at the top level and is not supposed to be embedded in another FD. All the paths appearing within this FD are interpreted relative to the total FD. In logical terms, a total FD is a universe of reference for FUF and in programming language terms, a total FD is an environment.

### 14.1 FD Accessors

The following functions access the value of a path within a total FD.

```
(top-gdp fd path)

(top-gdpp fd path)

;; Example:
fd1 = ((a 1)
      (b ((fset (b1 b2))
          (b1 1)
          (b2 {^2 a}))))
      (c {^ b b2}))

(top-gdp fd1 {a}) --> 1
(top-gdpp fd1 {a}) --> (a 1)
(top-gdp fd1 {c}) --> 1                ;; follows indirections
(top-gdpp fd1 {c}) --> (a 1)
(top-gdp fd1 {b}) --> ((b1 1) (b2 {^2 a})) ;; can contain a path inside
(top-gdp fd1 {x y}) --> nil             ;; adds a null path in fd1

fd1 = ((a 1)
      (b ((b1 1) (b2 {^2 a})))
      (c {^ b b2}))
      (x ((y nil))))                ;; fd1 is modified

(empty-fd (top-gdp fd1 {x})) --> t                ;; but x is still unbound.

(top-gdp fd1 {a a1}) --> none    ;; cannot go below a leaf.
(top-gdp fd1 {c c1}) --> none    ;; idem
(top-gdp fd1 {b b3}) --> none    ;; forbidden by fset.
(top-gdpp fd1 {a a1}) --> none
```

Top-gdp - short for go-down-path - follows a path within a total FD and returns the FD pointed to by the path. It is guaranteed to return a proper FD, that is, it will never return a path to another location, but instead, it will follow indirections until a proper FD is found.

Top-gdpp is a companion function which instead of returning the FD pointed to by the path returns the pair containing the appropriate value. As for top-gdp, this function always follows indirections and returns a pair whose second element is a proper FD and never a path.

If requested the value of a feature which is not present (not yet instantiated) in the total FD, top-gdp (or top-gdpp) will construct the path up to the requested level as shown in the example with the request for the value f the path x y which adds the path (x ((y nil))) into the total fd. This is the manifestation of the fact that the value nil adds no information to an FD. The path (x ((y nil))) does not change the meaning of the FD since it adds no instantiation.

The function filter-nils is used to remove these useless empty fds from an fd, and return only the useful information stored in an FD.

```
(filter-nils '((a nil) (b ((c 1) (d nil))) (x ((y nil)))))
=>
((b ((c 1))))
```

If the path requested is not compatible with the grammar, the value **none** is returned. This can happen in 2 cases: when the requested path extends below a leaf of the FD (for example the path a a1 above) or when the requested path does not satisfy an FSET declarations (for example the path b b3 above).

Top-gdpp can either return a pair from the total FD or none. When the result is a pair, it is possible to perform a (setf (second pair) new-value) and the total FD will be appropriately modified - as long as the new value is a valid FD. This is one way to “patch” a total FD with new values.

## 14.2 FD Relocation

Note that the value returned by top-gdp is not a valid total fd in itself because it can contain paths that point within the total FD. This is shown above in the result of (top-gdp fd1 b) which contains a path ((b1 1) (b2 {^2 a})). If you are interested in observing a sub-fd as a self-contained entity with no pointers into an external environment, you must use relocate instead of top-gdp. Relocate extracts a sub-fd from a total FD and turns it into a stand-alone total FD.

```
(relocate totalfd path)

Example:
(relocate '((a ((a1 {^ a2})
                (a2 ((x 2)))
                (a3 {a a1})
                (a4 {b})
                (a5 {^2 c})))
          (b {a a1})
          (c ((c1 1)))))
=>
((a1 {^ a2})      <---- NOTE keep relative path
 (a2 ((x 2)))
 (a3 {a1})        <---- NOTE updated path
 (a4 ((x 2)))     <---- NOTE loose conflation a4/a2 because went out of a
                  scope.
 (a5 ((c1 1))))  <---- NOTE resolved path
```

This function performs the following adjustments:

- Extract the sub-fd with top-gdp.
- Analyzes all paths found in the sub-fd. Absolute paths which point within the sub-fd are updated to make them relative to the new root. For example, above, the path (a3 a a1) is update to (a3 a1).
- Paths which point outside the scope of the new root are resolved and replaced by the value pointed to. For example, (a5 {^2 c}) is replaced by (a5 ((c1 1))).



- Relative paths which point inside the total FD remain relative paths. For example,  $(a1 \{^ a2\})$  remains relative.
- Conflations which are established indirectly through a feature which lies outside the scope of the sub-fd are lost. This is the case above between  $a1$  and  $a4$ . In the original total FD,  $a1$  and  $a4$  are unified, because  $\{a a4\}$  points to  $\{b\}$  which points to  $\{a a1\}$ . During the extraction analysis, it is found that  $a4$  points outside the scope of  $a$  (to  $b$ ), and therefore, the path is resolved. In this case, information is lost, since in the result,  $a1$  and  $a4$  happen to have the same value, but they are not conflated.

The reverse of relocate is insert-fd, which inserts a total FD within a larger total FD. The function also performs analysis of all the paths found within the inserted FD to adjust them to their new environment.

```
(insert-fd fd total subfd-path)

;; Insert Fd into TOTAL at location subfd-path.

;; Example:
(insert-fd '((a {b}) (b 1) (c {^ b}))
           '((b 2) (c ((x 1))))
           {c})
;; =>
((b 2)
 (c ((x1 1)
      (a {c b})
      (b 1)
      (c {c b}))))
;; <----- NOTE the inserted FD is unified with the
;;           existing sub-fd
;; <----- NOTE updated path.
;; <----- NOTE relative path is resolved
```

The main features of this function are illustrated by the example:

- The new FD is unified with the existing sub-fd at location subfd-path in the total FD (so the FD  $((x1 1))$  is enriched with the new FD).
- Paths in the inserted FD are adjusted to their new environment. Absolute paths are prefixed by subfd-path, and relative paths are resolved and made absolute relative to the new root (e.g.,  $(c \{^ b\})$  is resolved to  $(c \{c b\})$ ).

## 14.3 FD Normalization

The FD accessors and relocation functions, top-gdp, insert-fd and relocate, only work when the FDs are in normal form. Normal form is intuitively the “shortest notation” of an FD. For example, the normal form of  $((a ((a1 1))) (a ((a2 2))) (\{a a3\} 3))$  is  $((a ((a1 1) (a2 2) (a3 3))))$ . A normal form also does not contain disjunctions. The intention is to define a normal form for input FDs, with no extra nils, maximal factorization of features and no disjunctions.

The function normalize-fd takes as input an arbitrary FD and puts it in normal form. It is abstractly defined as:  $(\text{normalize-fd fd}) = (u \text{ nil fd})$ , since the unification function  $u$  returns an FD in normal form, and resolves disjunctions.

```
(normalize-fd fd)

;; Example:
(normalize-fd '((a ((a1 1)))
                (a ((alt (((a1 2))
                           ((a2 2))))))
                (a nil)
                ({a a3} 3)))
;; -->
(a ((a1 1)
     (a2 2)
     (a3 3)))
```

## 14.4 Lists of FDs

The functions `list-to-fd` and `top-fd-to-list` convert between lists of FDs and the FD encoding of lists discussed in Chapter 8.

```
(list-to-fd '( ((a1 1)) 2 ((x1 1) (x2 2)) ))
((car ((a1 1)))
 (cdr ((car 2)
       (cdr ((car ((x1 1) (x2 2)))
              (cdr none))))))

(top-fd-to-list '((car 1) (cdr ((car 2) (cdr none)))))
(1 2)

(top-fd-to-list (list-to-fd <l>))
<l>
```

# Chapter 15

## Reference Manual

For the sake of completeness, this chapter includes a list of all the functions, variables and switches that a user of FUF can manipulate. They are grouped under 8 categories. In each category, the list is sorted alphabetically:

1. Unification functions
2. Checking
3. Tracing
4. Complexity
5. Manipulation of the dictionary
6. Linearization and Morphology
7. Manipulation of FDs as Data-structures
8. Fine-tuning of the unifier

### 15.1 Unification functions

#### 15.1.1 `*lexical-categories*`

**Type:** variable

**Description:** The `*lexical-categories*` variable is a list of category names. These categories are those that are sent to the morphology component without being unified.

**Standard Value:** (verb noun adj prep conj relpro adv punctuation modal)

#### 15.1.2 `*u-grammar*`

**Type:** variable

**Description:** The `*u-grammar*` variable contains a Functional Unification Grammar. It is the default value to all the functions expecting a grammar as argument. It is a valid form if `grammar-p` accepts it.

#### 15.1.3 `*cat-attribute*`

**Type:** variable

**Standard value:** cat

**Description:** The `*cat-attribute*` variable contains a symbol. It is the default value for the `cat-attribute` argument to all the unification functions.

The CAT parameter is used to identify constituents in an fd when the `cset` attribute is not present. Through this mechanism, the unifier implements a breadth-first top-down traversal of the structures being generated.

By default, the `CAT` parameter is equal to the symbol `cat`. It is however possible to specify another value for this parameter. As a consequence, it is possible to traverse the same `fd` structure and to assign the role of constituents to different sub-structures by adjusting the value of this parameter. This feature is particularly useful when an `fd` is being processed through a pipe-line of grammars.

#### 15.1.4 `u`

**Type:** function

**Calling form:** (`u fd1 fd2 &optional limit success`)

**Arguments:** `p`

- `fd1` and `fd2` are arbitrary FDs. `fd1` cannot contain non-deterministic constructs, `fd2` can.
- `limit` is a number. The default value is 1000.
- `success` is a function of three arguments. It must be defined as: (`defun x (fd fail frame) ...`) where `fd` is an `fd`, `fail` is a continuation (a function) and `frame` is an object of type `frame`. The default value is the function `default-continuation`.

**Description:** `u` unifies `fd1` with `fd2` and passes 3 values to the `success` continuation: a resulting `fd`, a continuation to call if more results are needed and a “stack-frame” containing information needed to run the continuation. By default, `default-continuation` just returns the unified `fd`. `u` is a low-level function.

It is possible to limit the time the unifier will devote to a particular call. The `:limit` keyword available in all unification functions specifies the maximum number of backtracking points that can be allocated to a particular call. Using this feature it is possible to perform “fuzzy” unification. (Note that the appropriateness of a fuzzy or incomplete unification relies on the particular control strategy used of breadth-first top-down expansion.)

#### 15.1.5 `u-disjunctions`

**Type:** function

**Calling form:** (`u-disjunctions fd1 fd2 &key limit failure success`)

**Arguments:**

- `fd1` and `fd2` are arbitrary FDs. Both `fd1` and `fd2` can contain non-deterministic constructs, like `alt`, `ralt` and `opt`.
- `limit` is a number. The default value is 1000.
- `failure` is a function of one argument. It must be defined as: (`defun x (msg) ...`) where `msg` can be safely ignored.
- `success` is a function of three arguments. It must be defined as: (`defun x (fd fail frame) ...`) where `fd` is an `fd`, `fail` is a continuation (a function) and `frame` is an object of type `frame`. The default value is the function `default-continuation`.

**Description:** `u-disjunctions` unifies `fd1` with `fd2` and passes 3 values to the `success` continuation: a resulting `fd`, a continuation to call if more results are needed and a “stack-frame” containing information needed to run the continuation. By default, `default-continuation` just returns the unified `fd`. `u-disjunctions` is a low-level function. It is the only unification function accepting disjunctions in the input. For all other functions, if the input contains disjunctions, it should first be normalized by calling the function `normalize-fd`. The search for a unification works as follows: first one `fd` compatible with `fd1` is unified (as in `normalize`, in a blind search manner. Then, this `fd` is unified with `fd2`. If all tries with `fd2` fail, then the unifier backtracks and tries to find another `fd` compatible with `fd1`. Therefore, the choices in `fd1` are in general buried very deep in the search tree.

Refer to paragraph 15.1.4 for an explanation of the `limit` argument.

#### 15.1.6 `uni`

**Type:** function

**Calling form:** (`uni input-fd &key grammar non-interactive limit cat-attribute`)

**Arguments:**

- `input-fd` is an input fd. It must be recognized by `fd-p`. It must not contain disjunctions.
- `grammar` is a FUG. It must be recognized by `grammar-p`. By default, it is `*u-grammar*`.
- `non-interactive` is a flag. It is `nil` by default.
- `limit` is a number. It is 10000 by default.
- `cat-attribute` is a symbol. It has the value of `*cat-attribute*` by default.

**Description:** `uni` unifies *input-fd* with *grammar* and linearizes the resulting fd. It prints the result and some statistics if *non-interactive* is `nil`. It returns no value. *grammar* is always considered as indexed on the feature `cat` or of the `cat-attribute` argument. If *input-fd* contains no feature `cat` the unification fails. (cf. `unif` if this is the case.) If the input contains disjunctions, it should be normalized before being used (cf `normalize`).

Refer to paragraph 5.11 for an explanation of the `cat-attribute` argument. Refer to paragraph 15.1.4 for an explanation of the `limit` argument.

### 15.1.7 uni-string

**Type:** function

**Calling form:** `(uni-string input-fd &key grammar non-interactive limit cat-attribute)`

**Arguments:**

- `input-fd` is an input fd. It must be recognized by `fd-p`. It must not contain disjunctions.
- `grammar` is a FUG. It must be recognized by `grammar-p`. By default, it is `*u-grammar*`.
- `non-interactive` is a flag. It is `nil` by default.
- `limit` is a number. It is 10000 by default.
- `cat-attribute` is a symbol. It has the value of `*cat-attribute*` by default.

**Description:** `uni-string` works exactly like `uni` except that it returns the generated string as a lisp string and does not print it.

### 15.1.8 uni-fd

**Type:** function

**Calling form:** `(uni-fd input-fd &key grammar non-interactive limit cat-attribute)`

**Arguments:**

- `input-fd` is an input fd. It must be recognized by `fd-p`.
- `grammar` is a FUG. It must be recognized by `grammar-p`. By default, it is `*u-grammar*`.
- `non-interactive` is a flag. It is `nil` by default.

**Description:** `uni-fd` unifies *input-fd* with *grammar* and returns the resulting total fd. The result is determined. `uni-fd` prints the same statistics as `uni` if *non-interactive* is `nil`. *grammar* is always considered as indexed on the feature `cat-attribute`. If *input-fd* contains no feature `cat` the unification fails. (cf. `unif` if this is the case.)

Refer to paragraph 5.11 for an explanation of the `cat-attribute` argument. Refer to paragraph 15.1.4 for an explanation of the `limit` argument.

### 15.1.9 unif

**Type:** function

**Calling form:** `(unif input-fd &key grammar)`

**Arguments:**

- `input-fd` is an input fd. It must be recognized by `fd-p`.
- `grammar` is a FUG. It must be recognized by `grammar-p`. By default, it is `*u-grammar*`.

**Description:** `unif` unifies *input-fd* with *grammar* and returns the resulting total fd. The result is determined.

If *input-fd* contains no feature `cat`, `unif` tries all the categories returned by `list-cats` until one returns a successful unification.

`unif` checks *input-fd* with `fd-p` and it checks *grammar* with `grammar-p`.

#### 15.1.10 u-exhaust

**Type:** function

**Calling form:** (`u-exhaust fd1 fd2 &key test limit`)

**Arguments:**

- `fd1` and `fd2` are fds. They must be recognized by `fd-p`. `fd1` cannot contain disjunctions.
- `test` is a lisp expression. It is `T` by default.
- `limit` is a number. It is 10000 by default.

**Description:** Unifier exhaust. Takes 2 functional descriptions and returns the list of all possible unifications until `test` is satisfied. `Test` is a lisp expression which is evaluated after each possible unification is found. In `test`, refer to the list being built as `fug3::result`. This function does NOT recurse on sub-constituents. It is at the same level as `u` or `u-disjunctions` (low-level function). If `test` is `nil`, this function will generate all possible unifications of `fd1` and `fd2`.

Refer to paragraph 15.1.4 for an explanation of the `limit` argument.

#### 15.1.11 u-exhaust-top

**Type:** function

**Calling form:** (`u-exhaust-top input &key grammar non-interactive test limit`)

**Arguments:**

- `input` is an fd. It must be recognized by `fd-p`. It cannot contain disjunctions.
- `grammar` is a FUG. It must be recognized by `grammar-p`. By default it is `*u-grammar*`.
- `non-interactive` is a flag. If it is `nil`, statistics are printed after each unification. Otherwise the function works silently. The default value is `nil`.
- `test` is a lisp expression. It is `T` by default.
- `limit` is a number. It is 10000 by default.

**Description:** Unifier exhaust with recursion. This function works like `u-exhaust` except that it also recurses on the sub-constituents of the input. It keeps producing fds until `test` evaluates to a non-`nil`. `test` is evaluated in an environment where `fug3::result` is bound to the list of all fds found so far. If `test` is `nil`, this function will generate all possible unifications of `input` and `grammar`.

Refer to paragraph 15.1.4 for an explanation of the `limit` argument.

#### 15.1.12 uni-num

**Type:** function

**Calling form:** (`uni-num input n &key grammar limit`)

**Arguments:**

- `input` is an fd. It must be recognized by `fd-p`. It cannot contain disjunctions.
- `n` is an integer.
- `grammar` is a FUG. It must be recognized by `grammar-p`. By default it is `*u-grammar*`.
- `limit` is a number. It is 10000 by default.

**Description:** Unifies input with grammar and backtracks `n` times. Each time, the result is processed as per `uni`. Refer to paragraph 15.1.4 for an explanation of the `limit` argument.

## 15.2 Checking

### 15.2.1 fd-syntax

**Type:** function

**Calling form:** (fd-syntax &optional *fd* &key *print-warnings* *print-messages*)

**Arguments:**

- *fd* is a list of pairs. It is *\*u-grammar\** by default.
- *print-warnings* is a flag. It is *nil* by default.
- *print-messages* is a flag. It is *nil* by default.

**Description:** *fd-syntax* verifies that *fd* is a valid fd. If it is, it returns *T*. Otherwise, it prints helpful messages and returns *nil*. If *print-warnings* is non-*nil* it also print warnings for all the paths it encounters in the grammar. This is useful when you suspect that one path is invalid or pointing to a bad location. If *print-messages* is *nil*, no diagnostic messages are printed, and the function just returns *T* or *nil*. If it is non-*nil*, diagnostic messages are printed.

\*\*\*\*\*NOTE: The following table has not been updated for Version 5.2\*\*\*\*\*

"Diagnostics detected by <b>fd-syntax</b> "	
Message	Condition
Unknown type: <code>s</code> .	An fd must be either a legal leaf (symbol, number, string, character or array), a path or a valid list of pairs.
Unknow type: <code>s</code> . Should be a pair or a tracing flag.	Within a list of pairs fd, all elements must be either pairs or tracing flags.
<code>A</code> is a tracing flag. It cannot be used as an attribute in <code>A</code> .	Tracing flags are not legal attributes in a pair.
The attribute of a pair must be a symbol or a path: <code>s</code> .	Other types are forbidden.
An alt/ralt pair must have at most 4 args: (alt trace index demo (fd1 ... fdn)).	The disjunction construct being checked has more than 4 arguments.
An alt/ralt pair must have at most 4 args: (alt trace index demo (fd1 ... fdn)). There is no trace/index/demo flag in this pair.	The disjunction construct being checked has only 4 arguments, but one of the arguments is not properly formed.
An alt/ralt pair must have at most 4 args: (alt trace index demo (fd1 ... fdn)). There are no /index and demo/trace and demo/trace and index flags in this pair.	The disjunction construct being checked has only 3 arguments, but one of the arguments is not properly formed.
This alt/ralt pair must have one value: (alt (fd1 ... fdn)). There is no valid trace, index or demo flag in this pair.	No valid modifier is found in the current construct, and yet it has more than one value.
Value of alt/ralt must be a list of at least one branch: <code>s</code> .	A (alt ()) pair is invalid.
Value of special attribute ALT/RALT must be a list of valid FDs.	One of the fds in the branches of the pair is not valid.
An OPT pair must have at most 2 args: (OPT trace fd).	The OPT pair being checked has more than 2 arguments, as in (opt t a b).
This OPT pair must have at most 1 value: (OPT fd). (There is no valid tracing flag in this pair.)	A valid tracing flag must be either a symbol or a form (trace on symbol). The construct being checked has the form (opt f1 f2 ...) where f1 is not a valid tracing flag.
Value of OPT must be a valid FD: <code>s</code>	The fd part of the OPT is not a valid fd.
An FSET pair must be of the form (ATT VALUE): <code>s</code> .	The FSET pair has more than one value.
Value of special attribute FSET must be a list of atoms: <code>s</code>	One of the elements of the list is not a symbol.
A CSET pair must be of the form (ATT VALUE): <code>s</code>	The CSET pair has more than one value.
Value of special attribute CSET must be a list of symbols or valid paths: <code>s</code>	One of the elements of the list is not a symbol or a path.
A PATTERN pair must be of the form (ATT VALUE): <code>s</code>	The PATTERN pair can have only two values.
Value of special attribute PATTERN should be a list of paths or mergeable atoms.	<b>pattern</b> accepts a flat list of atoms, paths or mergeable constituents. A mergeable constituent is marked (* c).
A SPECIAL pair must be of the form (ATT VALUE): <code>s</code>	Special attributes can have only one value. <sup>1</sup>
An EXTERNAL pair must be of the form (ATT external-spec): <code>s</code>	External-spec can be simply the symbol <b>external</b> or a construct <b>#{external fctn}</b> .
— Warning: The argument of external must be a function.	In a <b>#{external fctn}</b> , fctn is not recognized as a defined function.
An UNDER specification must be an array of the form <b>#{UNDER symbol}</b> : <code>s</code>	An UNDER value has been found that does not have the proper form.
The argument of an UNDER specification must be a symbol: <code>s</code>	In a <b>#{under x}</b> , x must be a symbol defined in a type hierarchy.
— Warning: The argument of under does not have specializations defined. Use (define-feature-type <code>s</code> (spec1 ... specn)).	In a <b>#{under x}</b> , x must be a symbol defined in a type hierarchy. If the type hierarchy does not exist, it can be defined with the <b>define-feature-type</b> function.
A pair must be of the form (ATT VALUE)	The form being checked has more than two elements.
A value should be a valid fd.	In a pair (att value), value is not a valid fd.



## 15.2.2 fd-sem

**Type:** function

**Calling form:** (fd-sem &optional *fd grammar-p* &key *print-messages print-warnings*)

**Arguments:**

- *fd* is a syntactically valid fd. It must be recognized by *fd-p*. It is *\*u-grammar\** by default.
- *grammar-p* is a flag. It is T by default.
- *print-messages* is a flag. It is T by default.
- *print-warnings* is a flag. It is T by default.

**Description:** *fd-sem* verifies that *fd* is a semantically valid fd. If it is, it returns T. Otherwise, it prints helpful messages and returns nil. If *grammar-p* is non-nil *fd-sem* expects *fd* to be a grammar. It allows disjunctions in *fd*. In this case, *fd-sem* returns 3 values if *fd* is a valid grammar: T, the number of traced alternatives in the grammar, and the number of indexed alternatives.

If *grammar-p* is nil, *fd* is considered as an input fd. Disjunctions are not allowed. In any case, only one value is returned (T or nil).

\*\*\*\*\*NOTE: The following table has not been updated for Version 5.2\*\*\*\*\*

Diagnostics detected by <i>fd-sem</i>	
Message	Condition
— Warning: Disjunctions in input FD: <i>s</i>	<i>grammar-p</i> is nil and a disjunction has been found in <i>fd</i> .
— Warning: PATTERN or CSET should not be placed in input.	<i>grammar-p</i> is nil and a <i>pattern</i> or <i>cset</i> has been found in <i>fd</i> .
— Warning: ANY or GIVEN should not be placed in input.	<i>grammar-p</i> is nil and a <i>any</i> or <i>given</i> has been found in <i>fd</i> .
— Warning: EXTERNAL or UNDER should not be placed in input.	<i>grammar-p</i> is nil and a <i>external</i> or <i>under</i> has been found in <i>fd</i> .
Contradicting values for attribute <i>s</i> .	An attribute has been found with 2 different atomic values in the same branch of a disjunction. (for example, ((a 1) (a 2))).
This branch is contradictory: <i>s</i>	A branch in a disjunctive construct has been found with a contradictory value.

When *fd-sem* finds a problem in a grammar, it returns the path to the problem. Paths within grammars are different from paths in regular fds because of the presence of disjunctions. Paths have the same syntax, except that to go through an alt or ralt construct, additional information must be provided. The following syntax is used to denote the traversal of an fd with disjunctions:

PATH	:= { att-spec* }
ATT-SPEC	:= symbol   alt-spec   ralt-spec   opt-spec
ALT-SPEC	:= (alt <index> {<branch>})
RALT-SPEC	:= (ralt <index> {<branch>})
OPT-SPEC	:= (opt <index>)

Both index and branch must be numbers. The index information refers to the index of the alt, ralt or opt pair within its parent node. That is, given that a list of pairs can contain several alts at the same level, it is necessary to distinguish between them. The index information gives the position of the pair in the list, with index 0 referring to the first pair. If it is necessary to go further down an alt or ralt pair, then it is necessary to identify what branch must be followed. This information is given by the branch index. Not that a specifier (*alt index*) without a branch index MUST necessary be the last one in a path and refers to the whole alt pair. The function *alt-gdp* is used to traverse an fd with disjunctions using these extended paths as input. The function *get-error-pair* returns the first pair where *fd-syntax* would find an error.

### 15.2.3 fd-p

**Type:** function

**Calling form:** (fd-p *input-fd*)

**Arguments:**

- *input-fd* is an fd with no disjunctions.

**Description:** checks that *input-fd* is both syntactically and semantically a valid fd.

NOTE: Do not use fd-p on grammars.

### 15.2.4 grammar-p

**Type:** function

**Calling form:** (grammar-p &optional *fd print-messages print-warnings*)

**Arguments:**

- *fd* is a FUG. It is **\*u-grammar\*** by default.
- *print-messages* is a flag. It is **T** by default.
- *print-warnings* is a flag. It is **nil** by default.

**Description:** *grammar-p* verifies that *fd* is a valid grammar, both syntactically and semantically. If it is, it prints some statistics and returns **T**. Otherwise, it prints helpful messages and returns **nil**.

If *print-messages* is **nil** no statistics are printed.

If *print-warnings* is non-**nil** warnings are printed for all the paths encountered in the grammar. This is useful when you suspect that one path is invalid or pointing to a bad location.

NOTE: do not use *grammar-p* on input fds.

### 15.2.5 get-error-pair

**Type:** function.

**Calling form:** (get-error-pair *fd*)

**Arguments:**

- *fd* is an fd.

**Description:** *get-error-pair* checks the syntax of an fd. If the syntax is not correct, it returns the pair containing the first offending constituent of *fd*.

### 15.2.6 normalize-fd

**Type:** function.

**Calling form:** (normalize-fd *fd*)

**Arguments:**

- *fd* is an fd.

**Description:** *normalize-fd* prepares an fd that can contain disjunctions to be used as input to the standard unification procedures (that do not accept disjunctions, *i.e.*, all except **u-disjunctions**). If *fd* contains disjunctions, *normalize* will return one disjunction-free fd compatible with *fd*. It is best understood as basically an equivalent to the operation (**u nil fd**). If *fd* is not semantically correct (it contains contradictions), *normalize* will return **:fail**.

Normalize is also useful to put an fd in normal form with respect to the following constraint:

```
The fd ((a ((x 1))) (a ((y 2)))) is  
      ((a ((x 1)  
           (y 2))))  
in normal form.
```

## 15.3 Tracing

\*\*\*\*\*NOTE: This Section has not yet been updated for Version 5.2\*\*\*\*\*

\*\*\*\*\*Certain functions are missing\*\*\*\*\*

### 15.3.1 *\*all-trace-off\**

**Type:** variable.

**Description:** The *\*all-trace-off\** variable contains a flag that is recognized by the unifier and terminates the printing of all tracing messages. It must be placed in a valid position for a tracing flag.

**Standard Value:** %TRACE-OFF%

### 15.3.2 *\*all-trace-on\**

**Type:** variable.

**Description:** The *\*all-trace-on\** variable contains a flag that is recognized by the unifier and undoes the effect of the *\*all-trace-off\** flag, that is, it reenables all tracing messages. It must be placed in a valid position for a tracing flag.

**Standard Value:** %TRACE-ON%

### 15.3.3 *\*trace-determine\**

**Type:** variable.

**Description:** The *\*trace-determine\** is a switch enabling the printing of tracing messages on the determination stage. It indicates which TEST expressions are evaluated. When it is on and the determination stage fails in the context of a *uni* call, then the partially found sentence is linearized and printed. Cf *trace-determine* for the user-level interface to this variable.

**Standard Value:** T

### 15.3.4 *\*trace-marker\**

**Type:** variable.

**Description:** The *\*trace-marker\** variable contains a character. It is used to determine valid tracing flags: if the first character of the name of a symbol is *\*trace-marker\**, the symbol is a valid tracing-flag.

**Standard Value:** #%

### 15.3.5 *\*top\**

**Type:** variable.

**Description:** The *\*top\** variable is a switch enabling the printing of extensive debugging messages on the backtracking behavior of the unifier. Should be used for development only.

**Standard Value:** nil

### 15.3.6 *all-tracing-flags*

**Type:** function

**Calling form:** (*all-tracing-flags* &optional *grammar*)

**Arguments:**

- *grammar* is a FUG. It must be recognized by *grammar-p*. By default, it is *\*u-grammar\**.

**Description:** *all-tracing-flags* returns a list of all the tracing flags defined in *grammar*, in the order where they are defined in the grammar.

### 15.3.7 internal-trace-off

**Type:** function

**Calling form:** (internal-trace-off)

**Description:** `internal-trace-off` turns off the tracing of internal debugging information. Initially, no debugging information is printed.

### 15.3.8 internal-trace-on

**Type:** function

**Calling form:** (internal-trace-on)

**Description:** `internal-trace-on` turns on the tracing of internal debugging information. Initially, no debugging information is printed. Should be used for development only.

### 15.3.9 trace-disable

**Type:** function

**Calling form:** (trace-disable *flag*)

**Arguments:**

- *flag* is a tracing flag. A tracing flag must be an element of the result of `all-tracing-flags`.

**Description:** `trace-disable` disables the tracing flag *flag*. Initially, all tracing flags are enabled.

### 15.3.10 trace-disable-all

**Type:** function

**Calling form:** (trace-disable-all)

**Description:** `trace-disable-all` disables all tracing flags. Initially, all tracing flags are enabled.

### 15.3.11 trace-disable-match

**Type:** function

**Calling form:** (trace-disable-match *string*)

**Arguments:**

- *string* is a string.

**Description:** `trace-disable-match` disables all tracing flags whose names contain *string* as a substring. Initially, all tracing flags are enabled.

### 15.3.12 trace-enable

**Type:** function

**Calling form:** (trace-enable *flag*)

**Arguments:**

- *flag* is a tracing flag. A tracing flag must be an element of the result of `all-tracing-flags`.

**Description:** `trace-enable` enables the tracing flag *flag*. Initially, all tracing flags are enabled.

### 15.3.13 trace-enable-all

**Type:** function

**Calling form:** (trace-enable-all)

**Description:** `trace-enable-all` enables all tracing flags. Initially, all tracing flags are enabled.

### 15.3.14 trace-enable-match

**Type:** function

**Calling form:** (trace-enable-match *string*)

**Arguments:**

- *string* is a string.

**Description:** trace-enable-match enables all tracing flags whose names contain *string* as a substring. Initially, all tracing flags are enabled.

### 15.3.15 trace-off

**Type:** function

**Calling form:** (trace-off)

**Description:** trace-off turns off tracing. If no argument is provided, all tracing is turned off. Initially, tracing is off.

### 15.3.16 trace-on

**Type:** function

**Calling form:** (trace-on)

**Description:** trace-on turns on tracing.

Initially, tracing is off.

### 15.3.17 trace-determine

**Type:** function

**Calling form:** (trace-determine &key *on*)

**Description:** trace-determine turns on and off tracing for the determination stage.

When tracing is on for the determination stage, a message is printed indicating the location of the **any** found or the failed **test**. In addition, if the top-level function called is **uni**, the partially unified fd is linearized and printed to show the progression of the unifier.

### 15.3.18 trace-bk-class

**Type:** function

**Calling form:** (trace-bk-class &optional *on*)

**Description:** trace-bk-class turns on and off tracing of the special bk-class backtracking behavior.

When tracing is on for bk-class, a message is printed whenever a path of a bk-class category is caught by a backtracking point of the corresponding class. In addition, the number of frames that have been skipped thanks to the bk-class specification is printed. Examples of trace are provided in the section on bk-class.

### 15.3.19 trace-category

**Type:** function

**Calling form:** (trace-category *cat* &optional *on-off*)

**Arguments:** *cat* can be either a symbol (the name of a category), or the symbol **:all** specifying that all categories are to be traced, or a list of symbols (names of categories). If **on-off** is non-nil, the specified categories are traced, if it is nil, they are untraced.

**Description:** trace-category turns on and off tracing for all or certain categories.

When a category is traced, the unifier emits a message each time it starts unifying a constituent of that category. When **:all** is used, all categories are traced. This is particularly useful to identify in what constituent the unifier is failing and to follow the top-down breadth-first traversal of the constituents during unification. A useful sequence of calls is:

```
(trace-on)
(trace-disable-all)
(trace-category :all)
```

This sequence will trace the progression of the unifier at the level of the constituents, as in the following example:

```
> (uni a3)
>STARTING CAT DISCOURSE-SEGMENT AT LEVEL {}
>STARTING CAT UTTERANCE AT LEVEL {DIRECTIVE}
>STARTING CAT CLAUSE AT LEVEL {DIRECTIVE PC}
>STARTING CAT NP AT LEVEL {DIRECTIVE PC OBJECT}
>STARTING CAT NP AT LEVEL {DIRECTIVE PC IOBJECT}
>STARTING CAT NP AT LEVEL {DIRECTIVE PC SUBJECT}
>STARTING CAT LEXICAL-ENTRY AT LEVEL {DIRECTIVE PC VERB-CONCEPT}
>STARTING CAT NP AT LEVEL {DIRECTIVE PC OBJECT}
>STARTING CAT NP AT LEVEL {DIRECTIVE PC SUBJECT}
>STARTING CAT NP AT LEVEL {DIRECTIVE PC BENEF}
>STARTING CAT LEXICAL-ENTRY AT LEVEL {DIRECTIVE PC VERB-CONCEPT}
>STARTING CAT VERB-GROUP AT LEVEL {DIRECTIVE PC VERB}
>STARTING CAT PP AT LEVEL {DIRECTIVE PC DATIVE}
>STARTING CAT DET AT LEVEL {DIRECTIVE PC OBJECT DETERMINER}

{Used 108 backtracking points - 57 wrong branches - 103 undos}
John takes a book from Mary.
```

## 15.4 Complexity

### 15.4.1 avg-complexity

**Type:** function

**Calling form:** (avg-complexity &optional *grammar with-index rough-avg*)

**Arguments:**

- **grammar** is a grammar. It must be recognized by **grammar-p**. It is **\*u-grammar\*** by default.
- **with-index** is a flag. It is **T** by default.
- **rough-avg** is a flag. It is **nil** by default.

**Description:** **avg-complexity** computes a measure of the average complexity of a grammar. It tries to compute an "average" number of branches tried when the input to unification contains no constraint.

When *with-index* is **T**, all indexed alts are considered as single branches, when it is **nil**, they are considered as regular alts.

When *rough-avg* is **nil**, the average of an **alt** is the sum of the complexity of the first half of the branches. When it is **T**, the average is half the sum of the complexity of all branches.

### 15.4.2 complexity

**Type:** function

**Calling form:** (*complexity* &optional *grammar with-index*)

**Arguments:**

- *grammar* is a grammar. It must be recognized by `grammar-p`. It is `*u-grammar*` by default.
- *with-index* is a flag. It is T by default.

**Description:** *complexity* computes a measure of the complexity of a grammar. It tries to compute the worst case number of branches tried when the input to unification contains no constraint. The number it returns is equivalent to the number of branches the grammar would have in disjunctive normal form.

When *with-index* is T, all indexed alts are considered as single branches, when it is nil, they are considered as regular alts.

## 15.5 Manipulation of the Dictionary

### 15.5.1 \*dictionary\*

**Type:** variable

**Description:** The `*dictionary*` variable is a hash-table containing different types of entries. Each entry contains information on irregular morphological words.

The current dictionary contains entries for verbs, nouns and pronouns. It is defined in file LEXICON.L

The entries contain the following properties:

- verb : present-third-person-singular past present-participle past-participle
- noun : plural
- pronoun : subjective objective possessive reflexive.

### 15.5.2 lexfetch

**Type:** function

**Calling form:** (`lexfetch` *key property*)

**Arguments:**

- *key* is a non-inflected “root” form of a word. It must be a string.
- *property* is one of the properties defined in `*dictionary*` for the part-of-speech of the word.

**Description:** `lexfetch` fetches the inflected form of the word *key* from the hash-table `*dictionary*`. The properties accessible are those defined in `*dictionary*`.

### 15.5.3 lexstore

**Type:** function

**Calling form:** (`lexstore` *key property value*)

**Arguments:**

- *key* is a non-inflected “root” form of a word. It must be a string.
- *property* is one of the properties defined in `*dictionary*` for the part-of-speech of the word.
- *value* is the inflected form of *key* for *property*. It must be a string.

**Description:** `lexstore` stores the inflected form *value* of the word *key* in the hash-table `*dictionary*`. The properties accessible are those defined in `*dictionary*`.

## 15.6 Linearization and Morphology

### 15.6.1 call-linearizer

**Type:** function

**Calling form:** (`call-linearizer` *fd*)

**Arguments:**

- *fd* is a unified determined total fd. It must be accepted by `fd-p`.

**Description:** `call-linearizer` takes a complete determined fd in input and returns a string corresponding to the linearization of the fd.

### 15.6.2 gap

**Type:** feature.

**Description:** if a constituent contains the feature gap, it is not realized in the surface (it is a gap, still holding the place of an invisible constituent in the structure). It is used for implementing long-distance dependencies.

### 15.6.3 morphology-help

**Type:** function.

**Calling form:** (`morphology-help`)

**Description:** gives on-line help on what the morphology component can do.

## 15.7 Manipulation of FDs as Data-structures

### 15.7.1 fd-intersection

**Type:** function

**Calling form:** (`fd-intersection` *fd1* *fd2*)

**Arguments:**

- *fd1* and *fd2* are valid fds (recognized by `fd-p`). They represent lists as fds, using constituents `car` and `cdr`, and are terminated by a (`cdr none`).

**Description:** `fd-intersection` computes the intersection of two lists represented as fds, and returns the result as a regular lisp list.

### 15.7.2 fd-member

**Type:** function

**Calling form:** (`fd-member` *elt* *fdlist*)

**Arguments:**

- *elt* is any value acceptable as a value to an (attribute value) pair.
- *fdlist* is a valid fd (recognized by `fd-p`). It represents a list as an fd, using constituents `car` and `cdr`, and is terminated by a (`cdr none`).

**Description:** `fd-member` works as the lisp function `member` but on a list represented by an fd. It returns a list represented by an fd.



### 15.7.3 fd-to-list

**Type:** function

**Calling form:** (fd-to-list *fdlist*)

**Arguments:**

- *fdlist* is a valid fd (recognized by `fd-p`). It represents a list as an fd, using constituents `car` and `cdr`, and is terminated by a `(cdr none)`.

**Description:** `fd-to-list` converts a list from an fd representation to a lisp representation.

### 15.7.4 gdp

**Type:** function

**Calling form:** (gdp *fd path*)

**Arguments:**

- *fd* is a valid fd (recognized by `fd-p`).
- *path* is a valid path (that is a flat list of constituent names, starting with 0 or more `^`)

**Description:** `gdp` goes down the path *path* (hence its name: godownpath) and returns the fd found at the end of *path*. It is the only function that should be used to access sub-parts of an fd. `gdp` always returns a valid fd.

`gdp` works only if the special variable `*input*` is accessible and bound to the total fd containing *fd*. In normal environments, the function `top-gdp` should be used instead.

if *path* leads to a non-existent sub-fd, `gdp` returns:

- `none`: if the fd cannot be extended to include such a sub-fd (that's when we meet an atom on the way down)
- `any` : if the fd must be extended to include such a sub-fd (and exactly this sub-fd, that is only when the value is any)
- `nil` : otherwise (that is, an unrestricted fd).

### 15.7.5 gdpp

**Type:** function

**Calling form:** (gdpp *fd path frame*)

**Arguments:**

- *fd* is a valid fd (recognized by `fd-p`).
- *path* is a valid path (that is a flat list of constituent names, starting with 0 or more `^`)
- *frame* is a structure of type `frame`. By default it is `dummy-frame`, an empty frame.

**Description:** `gdpp` goes down the path *path* (hence its name: godownpathpair) and returns the pair whose value is the fd found at the end of *path*. It is the function that should be used to work as the basis to the `setf` of `gdp`, to set values to parts of an fd. `gdpp` always returns a pair whose second is a valid fd, and is never a path or `none` if *fd* cannot be extended to include *path*. (`gdpp *input* nil`) returns the pair `(*top* *input*)` (where `*input*` refers to the total fd).

`gdpp` works only if the special variable `*input*` is accessible and bound to the total fd containing *fd*. In normal environments, the function `top-gdpp` should be used instead.

if *path* leads to a non-existent sub-fd, `gdpp` extends (by physical modification) *fd* to include a path down to the required *path* if possible, or the function returns `none`. When the fd is modified physically, *frame* is updated (the field *undo*) to keep track of the modification.

### 15.7.6 top-gdp

**Type:** function

**Calling form:** (top-gdp *fd path*)

**Arguments:**

- *fd* is an fd.
- *path* is a valid path structure.

**Description:** works like **gdp** but considering *fd* as the total fd. Can therefore be used even if **\*input\*** is not bound in the current environment. This is the user-level function to access constituents of fds. .

### 15.7.7 top-gdpp

**Type:** function

**Calling form:** (top-gdpp *fd path*)

**Arguments:**

- *fd* is an fd.
- *path* is a valid path structure.

**Description:** works like **gdpp** but considering *fd* as the total fd. Can therefore be used even if **\*input\*** is not bound in the current environment. This is the user-level function to access pairs in fds. .

### 15.7.8 alt-gdp

**Type:** function

**Calling form:** (alt-gdp *fd path*)

**Arguments:**

- *fd* is an fd.
- *path* is a valid path structure with disjunction specifiers.

**Description:** works like **gdp** but works even if *fd* contains disjunctions (**alt**, **ralt** or **opt** constructs). paths when disjunctions are allowed are different from paths in regular fds. They have the same syntax, except that to go through an **alt** or **ralt** construct, additional information must be provided. The following syntax is used to denote the traversal of an fd with disjunctions:

```
path      := { att-spec* }
att-spec  := symbol | alt-spec | ralt-spec | opt-spec
alt-spec  := (alt <index> {<branch>})
ralt-spec := (ralt <index> {<branch>})
opt-spec  := (opt <index>)
```

Both index and branch must be numbers. The index information refers to the index of the **alt**, **ralt** or **opt** pair within its parent node. That is, given that a list of pairs can contain several alts at the same level, it is necessary to distinguish between them. The index information gives the position of the pair in the list, with index 0 referring to the first pair. If it is necessary to go further down an alt or ralt pair, then it is necessary to identify what branch must be followed. This information is given by the branch index. Not that a specifier (**alt index**) without a branch index must necessary be the last one in a path and refers to the whole **alt** pair.

This function is less robust than **gdp**, it does not check for cycles, does not check for fsets violations and does not check for user-defined special types.

### 15.7.9 list-to-fd

**Type:** function

**Calling form:** (list-to-fd *list*)

**Arguments:**

- *list* is a regular lisp list.

**Description:** list-to-fd converts a list from a lisp representation to an fd representation.

Note: the notation l 3 a refers to the third element of the list l if l is a list in fd form. That is, the path l 3 a is equivalent to l cdr cdr car a.

## 15.8 fine tuning of the unifier

### 15.8.1 \*any-at-unification\*

**Type:** variable

**Description:** if \*any-at-unification\* is nil, and the unifier encounters a pair (attribute **any**) in the grammar, and no feature **attribute** exists in the input, the unification succeeds and the input is enriched with the pair (attribute **any**). Only at the determination stage, it is checked whether **anys** remain in the total fd. If it is the case, the unification fails, and the unifier backtracks.

if \*any-at-unification\* is non-nil, the test to decide whether the feature attribute exists or not is performed immediately on the non-determined fd. The result may be incorrect, but it is much faster. The result is assured to be correct if the feature tested is one that is never instantiated by the grammar, and is expected to be provided in the input.

**standard value:** t

### 15.8.2 \*use-given\*

**Type:** variable

**Description:** if \*use-given\* is t, **given** and **under** constructs have their usual semantics. If it is nil, then **given** is considered as a normal symbol and a construct **#(under s)** is considered equivalent to simply **s**.

**standard value:** t

### 15.8.3 \*use-any\*

**Type:** variable

**Description:** if \*use-any\* is t, **any** constructs have their usual semantics (as determined by \*any-at-unification\*). If it is nil, then **any** is considered as a normal symbol.

**standard value:** t

### 15.8.4 \*keep-cset\*

**Type:** variable

**Description:** if \*keep-cset\* is nil, the determination stage removes all the **cset** features from the total fd. If it is t it keeps them.

**standard value:** nil

### 15.8.5 \*keep-none\*

**Type:** variable

**Description:** if \*keep-none\* is nil, the determination stage removes all the pairs whose value is **none** from the total fd. If it is t it keeps them.

**standard value:** t

### 15.8.6 \*agenda-policy\*

**Type:** variable

**Description:** can be either `:force` or `:keep`. Determines what to do with frozen alts at the end of unification:

- if `:keep`: keep them unevaluated in result
- if `:force`: force their evaluation at determination time.

the value should only be `:keep` in situations where several grammars are used in a pipeline architecture, and the frozen decisions are passed along from one stage to the next.

**standard value:** `:force`

# Appendices

# Appendix A

## Installation of the Package

### A.1 Finding the Files

You need to find out on which machine and under which directory the system is available. You also need to know how to run common lisp on that machine.

The file `fug5.l` will load all the required modules. examples are in the files `gr0`, `gr1` and up for the grammars, and in files `ir0`, `ir1`, ... and up for the inputs. the examples are of increasing complexity.

```
;; to try the examples, type:
```

```
lisp> (load "gr0")
t
lisp> (load "ir0")
t
```

### A.2 Porting to a New Machine

The program is contained in 29 files of source and 37 files of examples. All the source files should be grouped in a directory, that we will call here "`$fug5`", and the example files in a subdirectory of "`$fug5`" called "examples".

Once this is done, you probably need to edit the file `fug5.l`. This file loads all the required modules and defines a few functions useful for compiling or loading the package. In the file `fug5.l`, the function `require` is used to load all submodules. `require` takes as first argument the name of a module, and accepts a second optional argument, the name of the file containing that module.<sup>1</sup>

You must change the second arguments of all the `require` statements in file `fug5.l` and update there the name of the directory, from `$fug5` to the name of your directory.

You also need to edit the first line of the functions `compile-fug5` and `reload-fug5` and change there the name of the directory from `$fug5` to the new name.

When the file `fug5.l` is updated, load it in your common-lisp environment and follow these 4 steps:

```
(load "$fug5/fug5")
(in-package "FUG5")
(compile-fug5)
(reload-fug5)
```

Note to UNIX users: if you run Common Lisp under UNIX, and your version of Lisp can read environment variables and expands such variables in file names (for example, `(load "userx/file1")` is a valid statement, or `(load "$var/file2")`), then you don't need to edit the file `fug5.l`. All you need to do is to define the

---

<sup>1</sup>on Common Lisp II, the function `require` is generally not supported anymore. Most of the commercial implementations of Common Lisp, however, still have a form of `require` available.

environment variable "fug5" to be the complete pathname of the directory containing the source files. If you use Lucid Common Lisp or Franz Inc's Allegro Common Lisp, init files are provided in the system to redefine load so that it can recognize shell variables. These files are named lisp-init.l and lisp-init-acl.l respectively.

Once this installation is done, all you need to do to load the package is (load "\$fug5/fug5") (with \$fug5/ replaced by the name of your directory if you are not under unix).

## A.3 Packages

\*\*\*\*\*NOTE: This section has not been updated for Version 5.2 \*\*\*\*\*

\*\*\*\*\*Precise information can be obtained by inspecting file fug5.l in the distribution\*\*\*\*\*

The whole package is loaded in package 'fug5. The easiest way to access it is to type:

```
(in-package "FUG5")    ;; note the upper-case
or
(use-package "FUG5")
```

The following symbols are exported from package 'fug5 (they are the external symbols of the package, cf Steele 84, Chapter 11, pp171-192):

External symbols of package fug5	
File	Symbols
checker.l	fd-p fd-syntax fd-sem grammar-p get-error-pair
complexity.l	complexity avg-complexity
determine.l	*keep-cset* *keep-none*
external.l	*default-external-function*
graph.l	*any-at-unification* *use-given*
lexicon.l	*dictionary* lexfetch lexstore
linearize.l	call-linearizer morphology-help
path.l	gdp, gdpp top-gdp, top-gdpp alt-gdp
top.l	*u-grammar* *lexical-categories* *cat-attribute* uni uni-fd unif list-cats u u-disjunctions
trace.l	trace-on trace-off internal-trace-on internal-trace-off trace-category trace-enable trace-disable trace-enable-all trace-disable-all trace-enable-match trace-disable-match all-tracing-flags *trace-marker* *all-trace-off* *all-trace-on* *trace-determine* *top*
type.l	define-feature-type define-procedural-type reset-typed-features

In addition, the following symbols are external. These are the keywords used as names in the code:

External symbols of package fug5 (keywords)	
File	Symbols
top.l	=== * trace ^ alt any cat control cset dots *done* gap given index lex mergeable none opt pattern pound punctuation test

All these symbols are documented for reference in section 15. If you use the package fug5 in another package, only these symbols will be imported.



# Appendix B

## Advanced Features

### B.1 Advanced Uses of Patterns

In addition to constraining the ordering of constituents, the pattern unifier can be used to enforce the unification of constituents. The classical example is given by the **focus** constituent. There is good linguistic evidence that the focus of a sentence tends to occur first in a sentence. To represent this constraint, a grammar can include the following directive:

```
(pattern (focus dots))
```

That is, a sentence should start with its focus. Now, we also know that a sentence at the active voice should start with its subject, that is its **prot** constituent. this is expressed by:

```
(pattern (prot ... verb ...))
```

If both constraints are to be satisfied, we need to say that **focus** and **prot** are actually the same constituent, otherwise, the 2 patterns are incompatible. That is, the constituents **focus** and **prot** need to be unified. This mechanism would be quite expensive to implement for all constituents, and would need to meaningless attempts most of the time. Therefore, to allow this kind of unification to occur, the current unifier requires the pattern to include a special directive, indicating that a constituent can be unified with other constituents to make two patterns compatible. The notation used is: **(\* constituent)**.

```
;; Example:  
(pattern ((* focus) dots))  
(pattern (prot dots verb dots))
```

are compatible, and require the unification of the constituents **focus** and **prot**. Note that **prot** needs not be “stared” to be unified with **focus**. The notation can be understood as specifying that **focus** is a kind of “meta-constituent”.

Another advanced usage of mergeable constituents in patterns, is to infer a constraint on a constituent from its location in the pattern. A simple example is given here to enforce the decision to capitalize the first constituent of a sentence (the feature **punctuation** is interpreted by the linearizer as explained in Section 11.6):

```
((pattern ((* first-constituent) dots))  
 (first-constituent ((punctuation ((capitalize yes))))))
```

### B.2 Advanced uses of CSET

Note that CSET is rarely used in small to medium grammars, and most often used when you DO NOT want a sub-fd to be unified as a constituent, even though it is mentioned in a pattern or it contains a feature (cat xx).

When a CSET feature is specified, the order of the constituents can be important to make unification more efficient. The unifier traverses the input fd breadth-first identifying constituents at each level. Within the same level, the CSET feature when present specifies in which order the constituents must be unified. Therefore,

if there is a constituent known to be easy to unify, and whose value condition the unification of the brother constituents, it should be unified first, and placed first in the CSET. This way, the CSET feature can be used to optimize the work of the unifier.

```
((cat hard)
 (a \#)      ; is hard to unify
 (b \#)      ; is hard to unify
 (c \#)      ; is easy to unify and constrains the unification of a and b
 (cset (c a b))) ; unify c first, then a and b.
```

Explicit CSET traversal becomes very important when designing a large grammar such as SURGE. The incremental CSET specification discussed in Section 5.7 become then particularly important.

## B.3 Long Distance Dependencies and the GAP feature

The special feature **gap** is used to indicate that a constituent must not be realized in the surface text. If a constituent contains an attribute **gap** with any non-NONE value, the linearizer will skip it.

This device is used to implement long-distance dependencies in grammars. For example, in a relative clause, the relative pronoun can be viewed as the marker of the relativization, and the relative clause as a complete clause, with one constituent elided. Thus, in *The man whom I know*, the relative clause would have the structure *I know the man* and the constituent *the man* would be a **gap**, whereas the relative pronoun *whom* would inherit its properties.

## B.4 Specifying Complex Constraints: the TEST and CONTROL Keywords

**test** and **control** are two “impure” specifications: they do not rely on the principle of unification to prevent a successful unification of 2 FDs. **control** should not be used except under extremely special circumstances. For the time being, it can be considered a synonym of **test**.

**test** is used to add a complex constraint on the result of a unification. A complex constraint refers to any Lisp predicate. If at the end of the unification the predicate is satisfied when applied to the resulting fd, the unification succeeds, otherwise it fails, and the unifier backtracks to find another solution.

The special character **'#**

**'** is used to refer to parts of the FD in the expression of the constraints. A **'#**

**'** must be followed by a valid path (either absolute or relative). The expression **#**

**(^ ^ a b)** is replaced by the value of the feature referred to by that path before the predicate is evaluated.

The order in which the **test** predicates will be evaluated is obviously not determined. Side effects are therefore **STRONGLY** discouraged within the body of the **test** constraints.

```
;; Examples:
((a 1)
 (test (equal \#\{a} \#\{b})))

;; is equivalent to the nicer:

((a 1)
 (b {a}))

((a 1)
 (test (numberp \#\{a})))
```

There is conceptually the same difference between **TEST** and **CONTROL** as there is between **ANY** and **GIVEN**: **TEST** constraints are tested at determination time, whereas **CONTROL** constraints are tested as soon as the unifier meets them. **CONTROL** is therefore in general much more efficient than **TEST**, but the results it provides are unpredictable in certain cases (if the features tested are given a different value later on during the unification, the result of the test could be different).

## B.5 Copying vs Conflation: The # notation

\*\*\*\*\*Jacques: here you go\*\*\*\*\*

## Appendix C

# Non Linguistic Applications of the Unifier: FUF as a Programming Language

Unification as used in the theory of functional unification grammars is a powerful mechanism that is not restricted to linguistic domains. It can be viewed as a “programming language” of its own. Actually, it is similar by many aspects to PROLOG. There are however some very specific features that make working with this version of unification well adapted to grammars, and not so well to more classic programming tasks.

### C.1 Dealing with Lists: The Member/Append Example

To make things clear, this implementation includes a “grammar” doing some list processing. The only operations presented are *member* and *append*. This grammar is in the directory *examples* in file GR5.L. It is printed here for easy reference for the discussion.

```
'((alt
  (((cat append)
    (alt append
      ;; First branch: append({},Y,Y).
      (((x none)
        (z {^ y})
        ;; This is to normalize the result of a (cat append):
        ;; it must contain the CAR and CDR of the result.
        (car {^ z car})
        (cdr {^ z cdr})))

      ;; Second branch: append({X/Xs},Y,{X/Z}):—append(Xs,Y,Z).
      ((alt (((x ((car any)))) ; this alt allows for partially
              ((x ((cdr any)))) ; defined lists X in input.
            ;; recursive call to append
            ;; with new arguments x, y and z.
            (cset (z))
            (z ((car {^ ^ x car})
                (cdr ((cat append)
                      (x {^ ^ ^ x cdr})
                      (y {^ ^ ^ y}))))))
            (car {^ z car})
            (cdr {^ z cdr}))))))

  (((cat member)
    (alt member
      (((x {^ y car}))
        ((y ((cdr any))))
        (m ((cat member)
            (x {^ ^ x})
            (y {^ ^ y cdr}))))))))))
```

This grammar is actually almost equivalent to the following PROLOG program:

```

member(X, {X|_}).
member(X, {_|Y}) :- member(X,Y)

append({}, Ys, Ys).
append({X|Xs}, Ys, {X|Zs}) :- append(Xs, Ys, Zs).

```

Note that the PROLOG form is much nicer! But there are reasons to look at the FUG version anyway. Here is how it works.

## C.2 Representing Lists as FDs

The first problem to handle lists with FUGs, is to represent lists as FDs, since FUGs can handle only FDs.

Quite simply, lists are represented as an FD with two features, CAR and CDR (with names ala Lisp).

```

;; The list (a b c) is represented by the FD:

((car a)
 (cdr ((car b)
       (cdr ((car c)
             (cdr none)))))))

;; The list (a (b c)) is represented by the FD:

((car a)
 (cdr ((car ((car b)
             (cdr ((car c)
                   (cdr none))))))
       (cdr none))))

```

### C.2.1 NIL and variables

Note in the previous example that the equivalent of the lisp atom NIL is NONE in the FD. NIL in an FD means “anything can come here” whereas NONE means “nothing can come here”. NIL therefore plays a role similar to uninstantiated variables in PROLOG.

The PROLOG expression {a X c} can be represented by the FD:

```

((car a)
 (cdr ((car nil)
       (cdr ((car c)
             (cdr none)))))))
<=>
((car a)
 (cdr ((cdr ((car c)
             (cdr none)))))))

```

The PROLOG expression {a b | Xs} can be represented by the FD:

```

((car a)
 (cdr ((car b))))

```

### C.2.2 The ~ notation

The car/cdr notation for lists is very awkward to use. The file FDLIST.L includes a mechanism to translate between the regular Lisp notation and the FD notation. It defines the macro-character “~” to indicate list values.

```

((cat member)
 (x a)
 (y ~(c b a)))
<=>
((cat member)
 (x a)
 (y ((car c)
      (cdr ((car b)
            (cdr ((car a)
                  (cdr none))))))))

```

Note that the “~” notation can be used only for completely specified lists. If some elements are uninstantiated, you must describe the list with the car/cdr notation.

The `~n` notation can be used to refer to elements of lists represented as FDs. The path `l 4` refers to the fourth element of the list `l`. It is equivalent to the path `l cdr cdr cdr car`.

### C.3 Environment and Variable Names vs. FD and Path

The notions of environment and variable in PROLOG or LISP correspond to the notion of “total FD” and path in Functional Unification. What we call a “total FD” is the highest level FD, the one corresponding to the path `()`. It is the FD corresponding to the input to the unifier, and that will be “determined” at the end of unification. This FD contains all the environment of a computation.

Variables are then just places or positions within this total FD.

If the total FD is the FD corresponding to `{a X c}`

```
((car a)
 (cdr ((cdr ((car c)
              (cdr none))))))
```

The variable `X` can be referred to by using the path `(cdr car)`

### C.4 Procedures vs. Categories, Arguments vs. Constituents

A program in FUG can be viewed as a collection of procedures, each procedure being represented by a category. In the `member` example of section III.1, an input containing the feature `(cat member)` will be sent to the `member` procedure.

Procedures expect arguments and return results. There is no notion of input and output in unification, as far as arguments are concerned. So we just consider arguments in general. For example, the `member` procedure has two arguments, called `X` and `Y` and represented in FUG notation by the constituents `X` and `Y` of the `(cat member)`.

The procedure `append` has three arguments, `X`, `Y` and `Z`. `Z` can be seen as the “result” of the procedure, or in functional notation: `Z = append(X,Y)`.

Note that, as in the corresponding PROLOG program, the FUG implementation of `member` and `append` is non-directional. All of the arguments can be partially specified, and the unification enforces the relation existing between them.

### C.5 The total FD Includes the Stack of all Computation

One problem with the way FUG work is that there is no notion of “environment” besides the total FD. Therefore, when a program works recursively, all the local variables that are normally stacked in an external environment are stacked within the total FD. At the end, the total FD contains the whole stack of the computation, and is pretty heavy to manipulate.

As an example here is the result of the simple call `append(a,b,c,d,Z)`:

```
((CAT APPEND)
 (X ((CAR A) (CDR ((CAR B) (CDR NONE))))))
 (Y ((CAR C) (CDR ((CAR D) (CDR NONE))))))
 (Z
  ((CAR A)
   (CDR
    ((CAT APPEND)
     (X ((CAR B) (CDR NONE)))
     (Y ((CAR C) (CDR ((CAR D) (CDR NONE))))))
   (Z
    ((CAR B)
     (CDR
      ((CAT APPEND)
       (X NONE)
       (Y ((CAR C) (CDR ((CAR D) (CDR NONE))))))
     (Z ((CAR C) (CDR ((CAR D) (CDR NONE))))))
    (CAR C))
```

```

(CDR ((CAR D) (CDR NONE))))))
(CAR B)
(CDR
  ((CAT APPEND)
   (X NONE)
   (Y ((CAR C) (CDR ((CAR D) (CDR NONE))))))
   (Z ((CAR C) (CDR ((CAR D) (CDR NONE))))))
   (CAR C)
   (CDR ((CAR D) (CDR NONE))))))
(CAR A)
(CDR
  ((CAT APPEND)
   (X ((CAR B) (CDR NONE)))
   (Y ((CAR C) (CDR ((CAR D) (CDR NONE))))))
   (Z
    ((CAR B)
     (CDR
      ((CAT APPEND)
       (X NONE)
       (Y ((CAR C) (CDR ((CAR D) (CDR NONE))))))
       (Z ((CAR C) (CDR ((CAR D) (CDR NONE))))))
       (CAR C)
       (CDR ((CAR D) (CDR NONE))))))
    (CAR B)
    (CDR
      ((CAT APPEND)
       (X NONE)
       (Y ((CAR C) (CDR ((CAR D) (CDR NONE))))))
       (Z ((CAR C) (CDR ((CAR D) (CDR NONE))))))
       (CAR C)
       (CDR ((CAR D) (CDR NONE))))))
  ))

```

Fortunately, the only thing of interest in this FD is probably the value of the constituents CAR and CDR of Z.

## C.6 Analogy with PROLOG programs

We have seen so far what aspects of FUGs are specific and different from other programming languages.

A program written using a FUG is very similar to a PROLOG program:

- The notion of success and failure in unification are equivalent to the “yes” and “no” of PROLOG programs.
- Simple statement can be combined using the connectives AND and OR: both FDs and PROLOG statements make use of conjunction and disjunction.
- Both notations rely heavily on unification, and refinement of partial descriptions to perform computations.

## C.7 Use of Set Values in Linguistic Applications

This discussion of FUGs as programming languages can appear frivolous. It is actually motivated by the desire to integrate more expressive features in linguistic grammars.

There are many different reasons to use set values in grammatical descriptions. For example, to describe a conjunction like “John, Mary and Frank” the set John, Mary, Frank appears as a good candidate. Many other applications for the category of set appear quite naturally when writing a grammar.

We want to be able to express grammatical constraints on such constructs within the framework of FUGs. We have found the procedures *member* and *append* to be quite useful in this attempt.

## Appendix D

# Non Standard Features of the Implementation and Restrictions

The current implementation includes features not available in other systems working with functional unification, and imposes restrictions. This section lists these non-standard aspects of the implementation. For each of the restriction, it is precised whether the checking functions (`fd-p`, `fd-sem` and `grammar-p`) detect the limitation or not.

### D.1 Typed Features

This implementation supports the definition of types of symbols as described in the section on typing.

### D.2 The FSET Special Attribute and Typed Constituents

This implementation supports the special attribute `FSET` to implement the notion of typed constituent and express completeness constraints. An fd `((fset (a b c)) (a 1))` can only have defined values (`non-none` values) for the attributes `a`, `b` and `c`. All other attributes are considered as having value `none`. Two `FSET` features can be unified, and the result is the intersection of the two values.

File `test1.1` in the example directory contains examples of use of `fset` feature.

### D.3 User-defined Types

User defined types with special unification procedures can be defined in this implementation. Refer to the `define-procedural-type` entry in the reference manual for details. File `test2.1` in the example directory contains examples of use of procedural types.

### D.4 Limits on Disjunction in Input

To work best, the unifier requires the input to be a simple FD, containing no disjunction (`alt`, `ralt` or `opt`). It can contain patterns. `tests` and `controls` are not allowed in input.

It is advised not to put `patterns`, `csets` or `anys` in the input fd. These constructs are indeed best viewed as devices used by the grammar to realize or enforce some constraints. The input should be left as “declarative” as possible, and therefore should not contain such constructs.

If disjunction are found in an FD given to `fd-sem`, a warning message is printed. `fd-sem` also issues warnings if its argument contains `patterns` or `csets`.

It is possible to use input fds containing disjunctions in two ways: one way is to first “normalize” the input fd and randomly choose one fd compatible with the disjunction-full fd but containing no disjunction at all.



The function `normalize-fd` performs this operation. The other way to use disjunctions in input is to use the low-level unifier function `u-disjunctions`. Note that in general, `u-disjunctions` can be extremely inefficient unless some control mechanisms are properly used.

## D.5 Mergeable Constituents in Patterns

An extension to the standard pattern unification mechanism is the use of “mergeable constituents”. A mergeable constituent in a pattern is noted `(* constituent-name)`. This notation indicates that when unifying the pattern containing it, this constituent can be “merged” or unified with another constituent that would need to be placed at the same position in the pattern.

For example, patterns `(a ... b)` and `(c ... b)` cannot be unified, because the first position of the unifying pattern would need to be both `a` and `c`. But patterns `((* a) ... b)` and `(c ... b)` can be unified, under the constraint that constituents `a` and `c` be unified (or “merged”). See also section 5.6 for a description of pattern unification.

## D.6 Indexing of Alternation

This implementation allows indexing of `alts`. The notation used is:

```
(alt trace-flag (index ... indexed-path) (branches+))
```

where each branch is a regular fd. The validity of the `indexed-path` is checked by the function `grammar-p`.

## D.7 Test and Control

It is possible to specify arbitrary constraints on the result of an unification within the grammar by using the constructs `test` and `control` described in section 4.7. The notation is:

```
(TEST lisp-expression)
```

where *lisp-expression* is an arbitrary lisp expression, where certain variables can be `#(path)`, and refer to the value of `(path)` in the determined result of the unification (see section 5.9 for a definition of the determination stage of unification).

Unification succeeds if the evaluation of *lisp-expression* in the environment of the determined result is non-`nil`. If it is `nil`, the unifier backtracks.

`control` works in a similar way, except that the *lisp-expression* is evaluated immediately when the unifier encounters the `control`, and therefore is evaluated in a non-determined fd.

Note that both `test` and `control` can be used only to enforce complex constraints but not to compute complex results to be added in the unification.

The function `grammar-p` does not check that the value of `test` and `control` is a valid lisp-expression.

## D.8 GIVEN and UNDER

The special value `given` and the related construct `under` are defined in this implementation. A feature `(att given)` is unified with an input fd, if the input contains a real value for attribute `att` at the beginning of the unification. When typed features are defined, a feature `(att #(under symb))` is unified with an input fd, if the input contains a value for attribute `att` which is *more specific* than `symb` at the beginning of the unification.

Note that if `symb` has no specializations defined in a type hierarchy, the notation `(att #(under symb))` is equivalent to `((att given) (att symb))`.

`given` and `under` are useful to check the presence of required features in inputs.

## D.9 EXTERNAL Specifications and Macros

The implementation supports the **EXTERNAL** construct which provides a form of delayed dynamic expression of constraints and a mechanism for a macro facility in grammars.

A feature (`att #(external fctn)`) in a grammar is interpreted as follows: the unifier calls the function `fctn` with one argument, the path at which the external feature is being unified. The function must return a valid fd `F`. The unifier then continue unification with this returned fd instead of the `as` if the feature had been (`a F`) in the first place.

It is therefore possible to dynamically decide, at run-time, what constraints will be used in the grammar.

## D.10 User-defined CAT Parameter

The **CAT** parameter is used to identify constituents in an fd when the `cset` attribute is not present. Through this mechanism, the unifier implements a breadth-first top-down traversal of the structures being generated.

By default, the **CAT** parameter is equal to the symbol `cat`. It is however possible to specify another value for this parameter. As a consequence, it is possible to traverse the same fd structure and to assign the role of constituents to different sub-structures by adjusting the value of this parameter. This feature is particularly useful when an fd is being processed through a pipe-line of grammars.

## D.11 Resource Limited Processing

It is possible to limit the time the unifier will devote to a particular call. The `:limit` keyword available in all unification functions specifies the maximum number of backtracking points that can be allocated to a particular call. Using this feature it is possible to perform “fuzzy” unification. (Note that the appropriateness of a fuzzy or incomplete unification relies on the particular control strategy used of breadth-first top-down expansion.)

## D.12 BK-CLASS Control Specification

The idea behind **bk-class** is to take advantage of the knowledge of where a failure occurs in the graph being unified to filter the stack of backtracking points. Note that a failure always happens at a leaf of the graph because of a conflict between two atomic values. The path leading to the point of the failure is called the “failure address”.

Note also that backtracking points are all associated with a disjunction construct in the grammar.

In FUF, you can annotate each disjunction to either catch (restart) or ignore (propagate higher on the stack) failures occurring at certain pre-declared failure addresses.

The details are:

1. Define classes of failure addresses as regular expressions on paths. These are called the bk-classes (for backtracking classes). All failures that occur at an address which is not a member of any bk-class are treated normally by all backtracking point.
2. Annotate certain disjunctions (and therefore the backtracking points they place on the stack) with the list of classes to which they are allowed to respond:

```
(alt (:bk-class c1 ... cn) ...)
```

3. During backtracking, a backtracking point checks if the failure address is part of a bk-class, and if it is, whether it is a member of the bk-classes it is declared to process. If it is, the b-point restarts the computation (by proceeding to the next alternative in the disjunction), otherwise it just ignores the failure and propagates it to the rest of the stack.

Through this mechanism, only the decisions that are “relevant” to the current failure need to be undone.

In practice, this mechanism when it is applicable provides huge performance gains, but it is rarely applicable: the problem is to define the bk-classes and it is a difficult task. You need to experiment a lot with it though as

in general, `bk-class` has lots of potential. The problem is deciding when it is applicable in the context of a big grammar.

Note that this mechanism relates to the “cut” of Prolog but it also has a big advantage: if the `bk-class` annotations are done consistently (which cannot be checked automatically), then the semantics of the grammar is not modified by the annotation in the sense that all the results produced by the “pure grammar” would also be produced by the grammar with the annotation.

Note also that this mechanism is implemented in such a way that when it is not used in the grammar, there is very little overhead in processing time.

An example of use of `bk-class` is provided in grammar GR7 .

## D.13 WAIT Control Specification

The idea is to freeze a decision until enough information is gathered by the rest of the unification process to make the decision efficiently. The annotation is:

```
(alt (:wait (<path1> <value1>) ...) ...)
```

When the unifier meets such a disjunction, it first checks if all the  $path_i$  have a value which is instantiated ( $value_i$  is used with typed symbols). If they are all instantiated, it proceeds as usual. If they are not, then the disjunction is put on hold (frozen) and the unification proceeds with the other constraints in the grammar. Whenever a toplevel disjunction is entered, the unifier checks if one of the frozen disjunctions can be thawed. If all the decisions to be made are frozen, any one is forced.<sup>1</sup>

`Wait` is probably the most powerful of all control constructs. The main benefit it provides is that it allows to order decisions dynamically based on what the input contains.

## D.14 IGNORE Control Specifications

The idea is to just ignore certain decisions when there is not enough information, there is already enough information or there are not enough resources. The 3 cases correspond to the annotations:

```
(alt (:ignore-when <path> ...) ...)
(alt (:ignore-unless <path> ...) ...)
(alt (:ignore-after <number>) ...)
```

**Ignore-when** is triggered when the paths are already instantiated. This is used when the input already contains information and the grammar does not have to re-derive it.

**Ignore-unless** is triggered when a path is not instantiated. This is used when the input does not contain enough information at all, and the grammar can just choose an arbitrary default.

**Ignore-after** is triggered after a certain number of backtracking points have been consumed. This indicates that the decision encoded by the disjunction is a detail refinement that is not necessary to the completion of the unification, but would just add to its appropriateness or value. **IGNORE-AFTER IS CURRENTLY NOT IMPLEMENTED.**

The main problem with these annotations is that their evaluation depends on the order in which evaluation proceeds, and that this order is not known to the grammar writer. But in conjunction with `wait`, this issue is not necessarily a problem as “wait” establishes constraints on when a decision is evaluated.

---

<sup>1</sup>NOTE: These 2 issues require further optimization: when do you awake a frozen disjunction and which one do you force first can be decided by some sort of dependency analysis, but one that can be done statically by building incrementally a complex data-structure, otw the cost of the analysis would not be justified. Practically, the first issue is most important: it determines how much overhead is added on ALL disjunctions when freeze is used. I am experimenting currently with a “granularity” control system, which controls how often the agenda of frozen disj. is checked. This is currently not implemented.

# Appendix E

## Changes from Version to Version

### E.1 New Features and Modifications in Version 3

- Paths are now represented with curly braces instead of simple lists.  
Old notation: (a (^ ^ b))  
New notation: (a {^ ^ b})  
The emacs routine found in file src/convert.el converts from the old format to the new format.
- Typed features are now allowed.
- Paths can now be in the left position of a pair: (a x) is legal.
- External is introduced.
- Procedural types are introduced.
- Fset is introduced.

### E.2 New Features and Modification in Version 4

- bk-class is introduced.
- cset and pattern can contain paths.
- ordinal and cardinal are added to the morphology.

### E.3 New Features and Modification in Version 5

- syntax of alt is transformed with keywords :demo, :trace, :wait, :bk-class, :ignore-when, :ignore-unless and :order acceptable in any order. Old syntax is still recognized.
- wait is introduced.
- ignore-unless and ignore-when are introduced.
- ^n and ~n notations are introduced
- The incremental cset features are introduced (cset (= a ...) (== b ...) (+ c ...) (- d ...)).
- The punctuation capitalize feature is recognized by the morphology/linearizer module.
- The def-alt, def-conj, def-grammar syntax is introduced with :& and :!.
- Def-test and test are introduced.
- Trace-level is introduced, plus other minor tracing functions.
- %break% feature is introduced.
- Digit (yes, no or roman) under ordinal and cardinal is now recognized by the morphology module.

- $\hat{n}^{\sim}$  notation is introduced.
- $\#$  notation is introduced (to copy paths instead of unify them).
- Functions `relocate` and `insert-fd` are defined.

# Bibliography

- [1] Kay M. Functional grammar. In *Proceedings of the 5th meeting of the Berkely Linguistics Society*. Berkeley Linguistics Society, 1979.
- [2] Karttunen L. Features and values. In *Proceedings of the 10th International Conference on Computational Linguistics (COLING 84)*, pages 28–33, Stanford, California, July 1984. ACL.
- [3] Knight K. Unification: a multidisciplinary survey. *Computing Surveys*, 21(1):93–124, March 1989.
- [4] Ait-Kaci H. *A Lattice-theoretic Approach to Computation Based on a Calculus of Partially Ordered Type Structures*. PhD thesis, University of Pennsylvania, 1984.
- [5] Winograd T. *Language as a Cognitive Process*. Addison-Wesley, Reading, Ma., 1983.
- [6] Kaplan R.M. and J. Bresnan. *Lexical-functional grammar: A formal system for grammatical representation*, pages 173–281. MIT Press, Cambridge, MA, 1982.
- [7] M.A.K Halliday. *An Introduction to Functional Grammar*. Edward Arnold, London, 1985.
- [8] Appelt D. E. Telegram: A grammar formalism for language planning. In *Proceedings of the Eighth National Conference on Artificial Intelligence*, pages 595 – 9, Karlsruhe, West Germany, August 1983.
- [9] Steele G.L.Jr de Kleer J., Doyle J. and Sussman G.J. *Explicit Control of Reasoning*, pages 93–116. MIT Press, 1979.

# Index

- %break%, 65
- & (notation), 28
- agenda-policy (variable), 98
- all-trace-off (variable), 89
- all-trace-on (variable), 89
- any-at-unification (variable), 97
- cat-attribute (variable), 26, 81
- default-external-value (variable), 45
- dictionary (variable), 93
- irreg-plurals (variable), 46, 48
- irreg-verbs (variable), 46, 48
- keep-cset (variable), 97
- keep-none (variable), 97
- lexical-categories (variable), 81
- top (variable), 89
- trace-determine (variable), 89
- trace-marker (variable), 67, 89
- u-grammar (variable), 81
- use-any (variable), 97
- use-given (variable), 97
- ^ notation, 95
- ^ notation (ambiguity), 17
- ^n notation, 15, 74
- ~ notation, 107
- ~n notation, 97, 108
  
- %break% , 74
  
- a-an (morphological feature), 46, 47
- absolute path, 15, 74
- adj, 46
- adv, 46
- agreement (subject/verb), 10, 12
- all-tracing-flags (function), 65, 69, 89
- alt (keyword), 10, 20, 50, 52, 66, 96
- alt (special attribute), 73
- alt-gdp (function), 87, 96
- ambiguity of the ^ notation, 17
- anonymous, 75
- any (special value), 25, 75, 95, 97, 104
- append, 106, 108
- arguments to procedures (in FUG as program), 108
- avg-complexity (function), 51, 92
  
- bk-class (annotation), 50, 52, 72, 76, 91
- branch (of an alt), 10
  
- call-linearizer (function), 94
- car (in FDs), 33, 107
  
- cardinal, 46
- case, 47
- cat (special attribute), 26, 46, 52, 75, 84
- cat attribute, 52
- categories-not-unified (function), 46
- category, 26
- cdr (in FDs), 33, 107
- clean-fd (function), 74
- clear-bk-class (function), 76
- common lisp, 100
- common noun, 10
- comparison prolog/FUG as program, 109
- compile-fug (function), 100
- complexity, 27, 51
- complexity (function), 51, 92
- con, 46
- conflation, 16
- conjunction, 15
- constituent, 10, 22
- constituent traversal, 10, 22, 25, 53, 60, 61, 71, 72, 76
- constraint (feature as), 15
- constraints (specifying complex), 104
- control (keyword), 76, 104, 111
- control in FUF, 72, 73
- control-demo (function), 70
- cset (keyword), 10, 22, 72, 103
- cset (special attribute), 40, 76
- cset (unification), 22, 24
  
- def-alt (construct), 27, 28
- def-conj (construct), 27, 28
- default (in alt), 10, 20
- define-feature-type (function), 41, 76
- define-procedural-type (function), 43, 76
- demo (annotation), 50, 70
- demonstrative, 47
- denotation (of FDs), 15
- desperation (when debugging FUF), 74, 75
- det, 46, 47
- determination, 25, 56, 61, 71, 83, 84, 97, 108, 111
- dictionary, 48
- digit (morphological feature), 46
- disjunction, 20
- disjunctive normal form, 93
- distance, 47
- DNF (disjunctive normal form), 51
- dots (in pattern), 21

- dots notation, 21
- draw-grammar (function), 29
- draw-types (function), 42, 76
- efficiency (of unification), 52
- EMACS (text editor), 75
- empty-fd (function), 77
- ending, 47
- environment (of a FUG as a program), 108
- equations, 15, 16
- examples, 9
- external, 45, 112
- external (keyword), 76
- external symbols, 101
- failure (of unification), 12, 20, 25, 67
- far, 47
- FD, 6, 9
- fd-intersection (function), 94
- fd-member (function), 94
- fd-p (function), 7, 9, 74, 88, 110
- fd-sem (function), 87, 110
- fd-syntax (function), 85
- fd-to-list (function), 95
- features, 15
- filter-nils (function), 78
- first (person), 47
- FSET (special attribute), 110
- fset (special attribute), 37, 39, 43, 75
- fuf-postscript (function), 29
- fug5 (package), 101
- fug5.l (file), 100
- gap, 94, 104
- gdp (function), 95
- gdpp (function), 95
- gender, 47
- get-error-pair (function), 88
- given (special value), 25, 42, 76, 104, 111
- gr0.l (file), 9, 100
- gr1.l (file), 100
- gr11.l (file), 14
- gr2.l (file), 100
- gr3.l (file), 100
- gr4.l (file), 14, 100
- gr5.l (file), 100
- gr7.l (file), 54, 113
- grammar organization, 26
- grammar-p (function), 7, 74, 88, 110
- graph (FD as), 17, 74
- hyper-trace-category (function), 65, 72
- ignore (annotation), 76
- ignore-unless (annotation), 50
- ignore-when (annotation), 50
- index (annotation), 50, 76
- indexing, 52, 111
- infinitive, 47
- initial failure, 63
- initialize-lexicon (function), 48
- insert-fd (function), 78
- instantiated features, 25, 107
- internal-trace-off (function), 90
- internal-trace-on (function), 90
- ir0.l (file), 100
- ir1.l (file), 100
- ir2.l (file), 100
- ir3.l (file), 100
- ir4.l (file), 100
- ir5.l (file), 100
- ir7 (file), 55, 56
- jumping (to a branch), 52
- leaf, 15
- lex (special attribute), 10
- lexfetch (function), 93
- lexical categories, 46
- lexicon.l (file), 48, 93
- lexstore (function), 93
- limit, 82
- limit (keyword), 76
- limits on disjunction in input, 110
- linearization, 6, 12, 21, 46, 94
- linearize.l (file), 46
- linearizer, 76
- list (as FDs), 107
- list-to-fd (function), 80, 97
- lists (as fds), 94, 95, 97
- member, 106, 108
- mergeable constituents (in pattern), 86, 103, 111
- modal, 46
- Morphology, 12, 49
- morphology, 46, 94
- morphology-help (function), 13, 47, 94
- near, 47
- nil (special value), 15, 78, 95, 107
- non-deterministic constructs, 22, 52
- non-interactive (keyword), 71
- non-standard features of implementation, 26, 110
- none (special value), 24, 38, 75, 78, 95
- normalize-fd (function), 16, 17, 79, 82, 88, 111
- noun, 46, 47
- number, 47
- objective, 47
- opt (keyword), 20, 52, 66, 96
- optional features, 20
- order (annotation), 20, 50
- order independence, 20
- ordering constraints, 10, 21
- ordinal, 46
- packages, 101
- pair (attribute/value), 15
- past, 47
- past-participle, 47
- path, 15, 95



path (flat description of FDs), 15  
 path (unification), 16  
 path-value (function), 69, 74  
 pattern (keyword), 10, 12, 21, 52, 86  
 pattern (special attribute), 40  
 pattern (unification), 21  
 person, 47  
 personal, 47  
 phrase, 46  
 plural, 47  
 porting to a new machine, 100  
 possessive, 47  
 pound (in pattern), 21  
 pound notation, 21, 76, 104, 111  
 prep, 46  
 present, 47  
 procedures (in FUG as program), 108  
 prolog, 106, 108  
 pronoun, 47  
 pronoun-type, 47  
 proper noun, 10  
 punctuation, 46–48  
  
 quantified, 47  
 question, 47  
  
 ralt (keyword), 20, 52, 96  
 recursion, 10, 22  
 reference to the FD in a test expression, 104  
 reflexive, 47  
 register-category-not-unified (function), 46  
 relative path, 15, 74  
 reload-fug5 (function), 100  
 relocate (function), 78  
 relpro, 46  
 require (lisp function), 100  
 reset-typed-features (function), 41, 76  
 restrictions, 26, 110  
 root, 47  
  
 search (through the grammar), 52  
 second (person), 47  
 set values in FDs, 109  
 set-path-value (function), 69  
 singular, 47  
 star notation in patterns, 86, 103  
 structure sharing, 16  
 sub-constituents, 10  
 subjective, 47  
 subsume (function), 41  
 sunder (special value), 41  
 syntax, 15  
  
 tense, 47  
 test (keyword), 25, 76, 104, 111  
 third (person), 47  
 top-fd-to-list (function), 80  
  
 top-gdp (function), 74, 77, 95, 96  
 top-gdpp (function), 77, 95, 96  
 total FD, 77  
 total fd, 25, 83, 84, 95, 97, 108  
 total-fd, 96  
 trace-alts (function), 65, 73, 75  
 trace-bk-class (function), 65, 72, 91  
 trace-bp (function), 65, 75  
 trace-category (function), 65, 72, 91  
 trace-cset (function), 65, 72  
 trace-determine (function), 65, 71, 91  
 trace-disable (function), 65, 69, 90  
 trace-disable-all (function), 65, 69, 75, 90  
 trace-disable-alt (function), 65, 69  
 trace-disable-match (function), 65, 69, 90  
 trace-enable (function), 65, 69, 90  
 trace-enable-all (function), 65, 69, 90  
 trace-enable-alt (function), 65, 69, 75  
 trace-enable-match (function), 65, 69, 91  
 trace-level (function), 65, 75  
 trace-off (function), 65, 75, 91  
 trace-on (function), 65, 75, 91  
 trace-wait (function), 65, 73  
 tracing, 63  
 tracing (local), 67  
 tracing (of alt), 66  
 tracing (of opt), 66  
 tracing flag, 67, 86, 89  
 tracing messages, 67  
 types in unification, 110  
 types-postscript (function), 42  
  
 u (function), 82  
 u-disjunctions (function), 82, 111  
 u-exhaust (function), 84  
 u-exhaust-top (function), 84  
 under (special value), 26, 41, 76, 111  
 uni (function), 7, 82  
 uni-fd (function), 7, 74, 83  
 uni-num (function), 84  
 uni-string (function), 83  
 unif (function), 7, 83  
 unification, 6  
 unification (overall mechanism), 10  
 unification functions, 81  
 Unknown category, 12, 46  
 unknown category, 76  
 use-package (function), 102  
 user-defined types, 110  
  
 value (morphological feature), 46  
 variables (in FUGs), 107, 108  
 verb, 46, 47  
  
 wait (annotation), 50, 52, 73, 76  
 wait (control annotation), 25