

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

ProGram - a Development Tool for GPSG Grammars

Roger Evans

ive Science Research Paper

no: CSRP 036

iversity of Sussex
ive Studies Programme
of Social Sciences

on BN1 9QN

ProGram - a Development Tool for GPSG Grammars*

=====

Roger Evans,
Cognitive Studies Programme,
University of Sussex,
Falmer, Brighton,
East Sussex.
February 1985

Abstract

The "ProGram" Grammar Development System is a computational tool designed for linguists attempting to develop grammars expressed in the "Generalised Phrase Structure Grammar" (GPSG) formalism. The GPSG theory is complex enough that ensuring that any but a small grammar behaves as expected is a difficult task. ProGram aims to help overcome this problem by providing a computational representation for GPSG grammars, and tools which allow the linguist to explore the effects of different parts of the grammar on the analytic structures assigned by it. This paper describes the design and implementation of ProGram, and gives examples of its use on a small GPSG grammar.

This paper describes the result of research supported by the Social Science Research Council (UK), grant number HR 7829/1 to the University of Sussex (principal investigator: Gerald Gazdar).

* This paper will also appear in "Linguistics", late 1985.

036

ProGram - a Development Tool for GPSG Grammars

=====

1. Introduction

Because natural languages are so rich, any serious attempt to provide realistic formal characterisation for even a small part of a language is faced with the problem that the necessary complexity of the formalism makes the theory difficult to test, evaluate and sometimes even express. By using computational aids in such cases, it should be possible to extend the linguist's manipulative power, so that the development of interesting complex theories becomes much easier.

This paper describes the "ProGram" grammar development system [Evans and Gazdar 1984] [1], an example of a computational tool designed to assist the development of grammars expressed in the "Generalised Phrase Structure Grammar" (GPSG) formalism of Gazdar et al. (see, for example, [Gazdar and Pullum 1982] [Gazdar et al. 1985]). GPSG aims to formally characterise the syntactic structure of natural languages [2], and has been developed to a point where a large number of interesting syntactic phenomena can be accommodated in an integrated formal framework. But, unsurprisingly, its extensive formal coverage is matched by its complexity. A typical [3] GPSG grammar has seven or eight distinct components which interact in a complicated fashion, providing a description of a language at a very high level. Although theoretical development can continue by focussing attention on only a small portion of the formalism at a time, confident prediction and validation of the complete syntactic coverage of any non-trivial GPSG grammar (and hence the general validity of the theory) is difficult and, if undertaken manually, prone to error. The ProGram system is designed to help overcome these problems, allowing more thorough validation of experimental grammars, and at the same time making it easier to write large grammars using existing GPSG techniques.

2. System Objectives

The principle design objective for the ProGram system was to produce a GPSG grammar development tool which was usable and useful. For the system to be usable, it needed to be accessible to linguists with little computational experience, easily portable from one machine or operating system to another and practical (in particular, in terms of speed and size). These points are largely implementation issues which will be discussed in more detail below. For the system to be useful, it needed to provide facilities which actually help a linguist attempting to develop a grammar and which are convenient to use. Such facilities can be loosely divided into two groups: facilities for specifying and modifying a grammar using ProGram, and facilities for examining and testing the grammar once it has been specified.

The first step in the design of the grammar specification language was to decide exactly which version of the GPSG formalism to use. Although there has always been widespread informal agreement about the general shape of the theory, tying the theory down to the formal detail which a computational tool requires in a manner which satisfied everyone concerned, could have posed a problem. However, the appearance of [Gazdar and Pullum 1982] (henceforth GP82), which gave a thorough, formally adequate (more or less) and largely uncontroversial account, provided an independent formulation of the theory suitable for use

a basis for ProGram.

The most recent account of GPSG, [Gazdar et al. 1985], differs from GP82 in various ways. Most of these are details of convention definitions, etc. which can be ignored for present purposes. The most significant change, however, is worth noting here. GP82 (and ProGram) makes use of tree-structure categories and features. These are more fully described in section 4 below. [Gazdar et al. 1985], categories are expressed as feature sets (rather than feature trees), allowing greater flexibility, but this facility is not provided in ProGram.

GP82 uses various notational conventions to describe different parts of a grammar and some of these, such as the notation for feature definitions, are well described in the paper. Others are not: the interpretation of the 'head' category marker, for example, is less well specified, while the abbreviation of category specifications and the implicit correspondence of categories in meta-rules are examples of very vague conventions which appeal to linguists but are informal and unusable. The ProGram system could simply have ignored all of these extra notations and assumed only the rigorous formal descriptions actually required. This would have led to a rule expressed informally as (2.1) (this is the GP82 example (32)) being expressed for ProGram as something like (2.2). [4]

(2.1) < 5; VP --> V,NP >

(2.2) [cat [bar 1] [head [major [verb +] [noun -]]]] -->
[cat [bar [lexical 5]] [head [major [verb +] [noun -]]]],
[cat [bar 2] [head [major [verb -] [noun +]]]]

Such an approach is clearly undesirable, and indeed hardly satisfies either the 'usefulness' criteria above (it is neither helpful nor convenient). The ideal case would allow the linguist to use all the notational devices, formal and informal, that are normally used. Grammar specification in ProGram was designed with this in mind; at every point the aim was to make a grammar look as much like the GP82 notation as possible. To achieve this, the informal conventions needed to be formalised as far as possible and abandoned otherwise. While the aim has not been perfectly achieved in all cases, the examples throughout this paper give some indication of the success, and the topic is discussed more fully below. Compare (2.3), a typical ProGram immediate dominance rule, with (2.1) and (2.2).

(2.3) 5: vp --> v,np.

Having designed a suitable specification language for grammars, it remained to decide what facilities for manipulating a grammar should be provided. What sorts of things might a linguist want to do using the system? The GPSG grammar is a high level specification of a language which is mediated through an ordinary context-free grammar [5] and one tends to think of a grammar in terms of the context-free rules it generates, rather than directly in terms of the strings it allows. Thus the system should provide at the very least some facilities for examining this underlying context-free grammar. A straightforward listing of all the rules would not be very informative as the grammars are large and complex; some form of interactive exploration would be more appropriate. The simplest approach would be to allow the linguist to ask whether particular context-free rules could be derived from the grammar as specified. The linguist could concentrate attention on likely problem areas of the grammar, where spurious rules might be generated, or expected rules omitted.

There are two main difficulties with this approach: firstly, it would be difficult and time consuming for the linguist to check out rules to the point where she could be confident the grammar contained exactly the expected rule, and secondly, given a 'correct' grammar, it would still be difficult to predict exactly what structures (if any) it assigns to given strings. This latter problem suggests that the system is bound to need some sort of parsing facilities, even though parsing is not its main function. And this suggests an alternative approach where the parser is used as the principle diagnostic tool rather than as a peripheral utility.

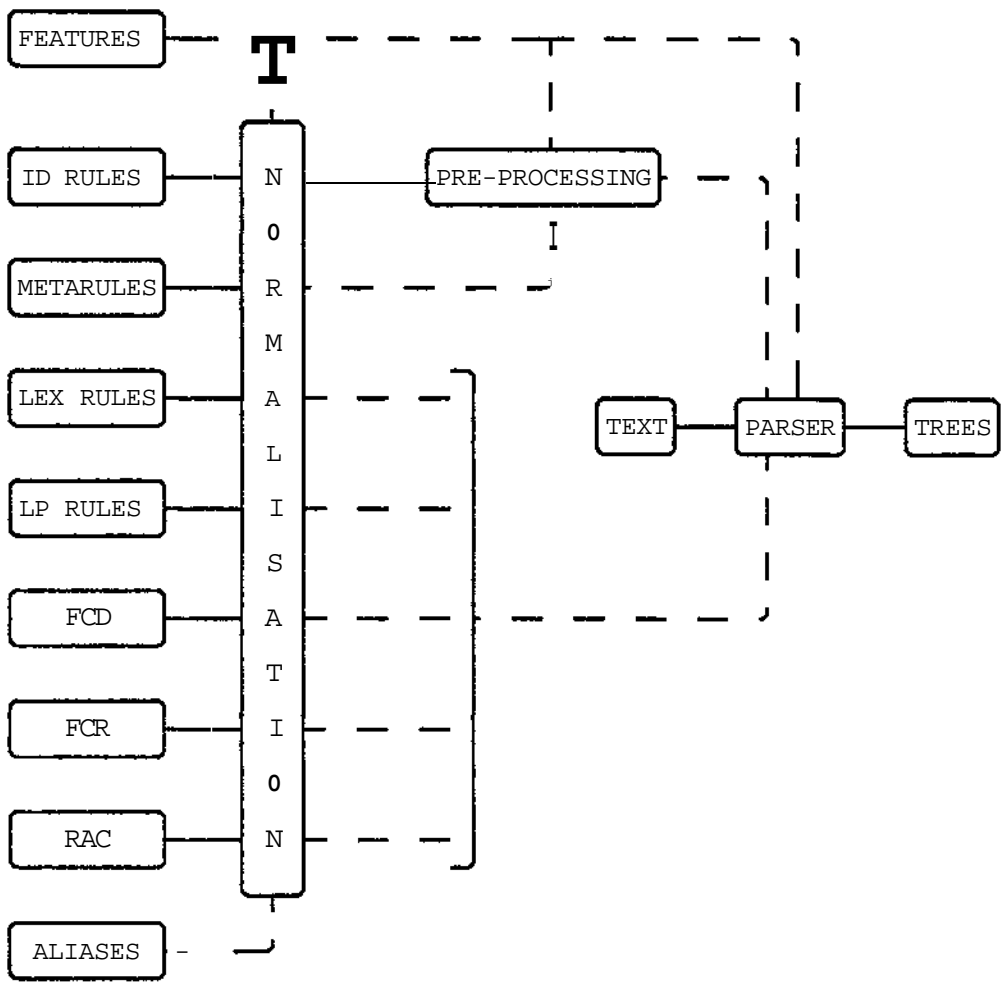
Suppose the system provided a parser which produced analyses of strings by generating context-free rules as needed from the GPSG data. Rather than ask about a context-free rule explicitly, the linguist attempts to parse a string whose analysis depends on the rule. A successful parse means that the rule is in the grammar. In fact, it means that several associated rules (all the others used in the analysis) are in the grammar too. Conversely, ungrammatical phrases could be used to ensure that certain rules were not in the grammar. More generally, the linguist would be able to test several different aspects of the grammar, and their interactions in assigning analyses, all at the same time, and the results of the tests would be available in a relevant context (the linguist would know not only that a rule is in the grammar, but also where in the analysis it can occur). Furthermore, grammar testing is achieved simply by providing test phrases, rather than making use of a complicated rule specification language.

Using a parser in this fashion is a common way of informally testing a grammar, and it often suffers from a lack of control and diagnostic information. Thus, for example, it is usually not possible to force the parser to attempt a particular analysis, or find out why an expected analysis did not occur. Where the context-free rules the parser is using are not immediately available to the linguist, but instead are being generated from higher-level data as the parser proceeds, such problems become more severe. But a special purpose parser, designed for the task of grammar development rather than analysis per se, can incorporate extra features to overcome these problems, resulting in a useful development tool. It is on this view of the problem that ProGram is based. The central component of the system is a parser which interprets a GPSG grammar and provides analyses of phrases, and which can be controlled by the user to allow more precise and detailed examination of chosen aspects of the grammar.

3. Overall Design

The following diagram shows the overall layout of ProGram. The boxes on the left hand side, and the one marked TEXT, are all data which the user must provide in some form. These are discussed in more detail in section 4. The box marked TREES contains the parse trees that are produced, and the remaining three boxes are the main functional units of the system.

OUTLINE STRUCTURE OF ProGram



Key: ————— Main data for process
 _ _ _ Background data for process
 (All data flows left to right)

The basic procedure for developing a grammar is as follows. Having created the initial data files for the grammar, the first task is to attempt normalisation. Most of the data modules have to be translated into an intermediate format for use with the parser. The normalisation process carries out translation, taking account of the built-in and user-specified notation conventions (see ALIASES in section 4), and checking that all the data is

are well-formed and, as far as possible, sensible. The user is informed of the problems encountered and the culprit files must be edited and re-compiled. Once all of these generally simple errors have been removed, the data is ready for use. Most of the data modules are used directly by the parser. The exceptions are the Immediate Dominance (ID) rules and ProGram carries out metarule expansion before parsing; the parser only uses the expanded set of ID rules. Furthermore, the Control Agreement Principle and the Head Feature Convention [6] are also compiled into the ID rules before parsing. These three processes are collected together in the PRE-PROCESSING block of the diagram. After normalisation then, the user may wish to expand the rules (or not, if there are no metarules, or they are being ignored for the present) and carry out the convention processing on them. The ID rules are then also be ready for the parser.

Although some minor errors will have been detected in these processes, the main diagnostic tool is the parser. Basic use of the parser is very simple: the LEX RULES module contains a lexicon, and the user enters phrases made up of words in this lexicon to be parsed. The parser can operate in three different modes, corresponding to different levels of control. In CONTROL mode, the parser runs completely independently, and produces one parse tree for the phrase. The user may set various tracing switches to monitor the parsing process, but has no control over it. MONITOR mode allows the user to ratify all the major decisions the parser makes. After a possible choice has been located, the user is informed and asked whether the parser should accept the choice. The user can agree to the choice or ask for an alternative choice (if any) to be sought. In this way, the actual calculation of possibilities is still automatic, but the user can direct the parser to search particular lines of search and so force particular analyses quickly. In INTERACTIVE mode the user specifies exactly what to do at each choice point, (for example, which ID rule to attempt to apply). Whenever a choice fails (that is, no rule is applied for some reason), the parser prints out diagnostic messages to the user exactly what went wrong (for example, the Foot Feature Principle is violated). Thus the user tells the parser exactly what to do, and may force it to follow paths it would not normally follow at all, and so can find out what parts of the theory is blocking expected analyses.

ProGram commands can be given with virtually no knowledge of the internal structure of all. For example, the command to parse phrases in a file 'sentences' and to print the trees in 'trees1' might be (3.1), while (3.2) is the command to expand a set of ID rules from a file called 'idrules', putting the expanded set into a file called 'idrules2', and using the file 'mrules' as the source of the (normalised) metarules.

(3.1) parse from sentences to trees1.

(3.2) expand from idrules to idrules2 using mrules for normalised metarules.

On the other hand, because this command structure is achieved using suitably defined Prolog operators, the more knowledgeable user can use ordinary Prolog commands, and easily incorporate new commands which modify the existing ones.

In addition to the utilities themselves, ProGram documentation is available in three forms. The main documentation takes the form of "The ProGram System" [Evans and Gazdar 1984], which gives detailed instructions for using the system. The system itself provides a quick help facility consisting

explanations of most of the commands and facilities in ProGram. Because this is built in to the standard system, it is always available regardless of local customisations etc. The more extensive 'help files' are also provided with the system, and can be read at the terminal using a suitable file-reading or editing package. In POPLOG, the screen editor, VED, is used to access these files in the same way as it is used for POPLOG help files. In other systems, it may be more convenient to have printed copies of them, or to use the manual which covers virtually the same material.

4. A ProGram Grammar

This section describes what a ProGram grammar actually looks like, in other words, what goes into each of the nine boxes down the left-hand side of the diagram above. The examples are taken from the example grammar which is used in section 5 and is repeated in full in the appendix. Although this grammar is quite simple, it incorporates ID/LP format, the Head Feature Convention, Control Agreement Principle, slash categories and metarules. This section will refer to a few features of the implementation language, Prolog. For full details see [Clocksin and Mellish 1981].

FEATURES

The most basic data for the grammar is the specification of the feature syntax, the FEATURES box. In GP82 a feature is a tree-structure with a root labelled with the feature name, and sub-trees which are other features. The feature data specifies for each feature its name and the names of permissible sub-features. The notation used by ProGram is virtually the same as in GP82, so that to specify a feature MINOR which takes a sub feature AGR, and either VFORM or CASE (but not both) one would put (4.1) in the FEATURES file.

```
(4.1) feature [minor agr {vform case}].
```

A FEATURES file will have many statements of this sort, defining the range of possible tree structures for features. Many of these features have special significance to ProGram (for example HEAD is used by the Head Feature Convention), so the overall structure of the feature system is usually similar to that given in GP82. The feature specification for this example grammar is given in (4.2).

```
(4.2) feature [root cat foot].
feature [cat bar head].
feature [bar {lexical 1 2}].
feature [head major minor].
feature [major {v n d p}].
feature [minor agr {vform pform nform}].
feature [agr {sing plur}].
feature [vform {fin pass}].
feature [pform {by to}].
feature [foot cat].
```

Additionally, the FEATURES data must specify which of these features is permissible as an actual syntactic category. (4.3) specifies that the feature ROOT is syntactic category, so that (4.4) is a valid category but (4.5) is not

(4.3) syncat root.

(4.4) [root [cat [bar 2] [head [major v]]] [foot]]

(4.5) [bar 2]

The feature specification gives ProGram the information required to do check categories in the other parts of the grammar. It also tells ProGram possible other features there can be in a category specification which is complete. Thus (4.7) may be just a more detailed instance of (4.6).

(4.6) [root [cat [bar 2] [head [major v]]]]

(4.7) [root [cat [bar 2] [head [major v] [minor [agr plur]]]
[foot [cat [bar 2] [head [major n]]]]]

ID RULES AND ALIASES

ID rules are specified in terms of syntactic categories (in the instances of the ROOT feature). The basic format of an ID rule is as in (4.8) and (4.9) is an example.

(4.8) <name>: <cat> -> <cat>, ... , <cat>.

(4.9) s: [root [cat [bar 2] [head [major v]]]] ->
[root [cat [bar 2] [head [major n]]]],
[root [cat [bar 3] [head [major v]]]]*

This is a traditional s -> np, vp rule, but expressed assuming the specification above. As was noted above, such rules are a little unwieldy. ProGram provides a mechanism to make them more manageable, through the use of aliases. An alias is an expression like (4.10).

(4.10) alias(s, [root [cat [bar 2] [head [major v]]]]).

This means, "if ever I write 's¹' as a category specification, I mean [root [cat [bar 2] [head [major v]]]]". If in addition there are aliases such as in (4.11) and (4.12), rule (4.9) can be written as (4.13).

(4.11) alias(vp, [root [cat [bar 1] [head [major v]]]]).

(4.12) alias(np, [root [cat [bar 2] [head [major n]]]]).

(4.13) s: s -> np, vp.

This has exactly the same effect as the original (note that the "s" before the colon is not affected by the alias - ProGram knows this is the name of the category).

More complex aliases can be written by using a functional notation (4.14). [7]

(4.14) alias(v(1), [root [cat [bar 1] [head [major v]]]]).
alias(v(2), [root [cat [bar 2] [head [major v]]]]).

Notice that these two aliases give alternative abbreviations for the

METARULES

The example grammar contains the three metarules shown in (4.19). To implement a simple passive, the "slash termination metarule" for introducing holes and "that-less relatives" (see [Gazdar et al. 1982]). Note the use of slash notation (A/B) which specifies category A with FOOT feature set to CAT of category B. This is an advanced but very convenient use of aliases discussed more fully in [Evans and Gazdar 1984].

```
(4.19) pass:      (VP1 --> ... , n(2) where VP1 is v(1))
                 ==>
                 (VP2 --> ... , opt(p(2,by))
                   where VP2 is v(1,pass),
                   VP1 matches VP2).

stm1:           (C1 --> C2, ...      where C1 is [root],
                 ==>
                 C2 is [root,[cat,[bar,2]
                 C1/C2 --> ... .

relcl:         (N1 --> ...          where N1 is n(1))
                 ==>
                 N1 --> ... , v(2)/n(2).
```

Each metarule has a name and consists of two nameless ID rules which are as above, except that the daughter lists may include "...", denoting a "multiset variable" which can match any number of arbitrary categories in a rule. In (4.19) there are further instances of Prolog variables. In STM1 RELCL variables are introduced simply to ensure that identical categories appear in the input and output ID rules to the metarule. Notice that C1 is defined as just [root] - that is, any category whatsoever. Similarly C2 is defined as just [root,[cat,[bar,2] - that is, any category whatsoever. In PASS, the variables are used to specify that categories "match". This implements the informal notion of correspondence between categories in the two rules. In this case the two mother categories are forced to be the same (even where they are not explicitly specified in the rule), except that the second one is passive. In general, a "matches" clause ensures that categories are identical except where they differ explicitly in the metarule.

LEXICAL RULES

The lexicon contains a mapping from categories to lexical items. A typical lexical rule looks like (4.20).

```
(4.20) vp2(fin): v0(fin) ->- sees.
```

This sets up the word SEES as an instance of a v0(fin) category. The rule is used for lexical subcategorisation: such a v0(fin) can only be introduced by a rule also labelled VP2 (above, the transitive vp rule). The extra (fin) in the rule name allows the user to distinguish this lexical rule from another (vp2(pass) for passive verbs) which may also be used in VP2. The full lexicon for this grammar is given in the appendix.

LINEAR PRECEDENCE RULES

The ID rules do not specify the orderings on the daughter categories. Instead this is done independently by the linear precedence (LP) rules. The example grammar contains the LP rules shown in (4.21).

```
(4.21) [root [cat [bar lexical]]] <<
        { [root [cat [bar 1]]]
          [root [cat [bar 2]]] } .
n(2) << v(1).
[root ~foot] << [root [foot cat]].
```

The use of curly brackets in (4.21) is a notational abbreviation to a group of mutually unordered categories to be treated like a single category.

LP rules need not specify complete categories, and an ordering of daughters is only valid if all the LP rules which match pairs of daughters are observed. Thus the first LP rule in (4.21) states that lexical categories precede non-lexical (that is, bar level 1 or 2) categories, the second noun-phrases precede verb-phrases and the third that slashed categories follow non-slashed categories (see [Farkas et al. 1983]). Note here that ~FOOT to mean "FOOT is unspecified".

FCD, FCR AND RAC RULES

The final three components of a ProGram grammar are Feature Co-occurrence Defaults, Feature Co-occurrence Restrictions and Root Admissibility Constraints. The first two of these are as described in GP82: FCD's specify default values for categories in ID rules which have to be met by valid instantiations. FCR's specify constraints on the structure of legal categories in general. The example grammar contains one FCD and three FCR's as shown in (4.22).

```
(4.22) fcd( vform, [foot], fin, free ).

        fcr( major, [foot], n, minor, [foot], nform ).
        fcr( major, [foot], p, minor, [foot], pform ).
        fcr( major, [foot], v, minor, [foot], vform ).
```

The FCD states that the default value for VFORM (but not the VFORM which is a sub-feature of FOOT), is FIN for lexical categories, and unconstrained for phrasal categories. The FCR's ensure that the MINOR features VFORM, NFORM only occur with appropriate MAJOR features. This excludes such illegal categories as passive nouns.

The RAC's specify additional featural constraints on root categories in any complete parse tree produced by ProGram. They allow the user to specify that certain parses are not interesting (for example, very underspecified categories), perhaps because they could never be part of a larger parse. In the simple example grammar there are no RAC's.

5. Examples using the grammar

This section provides examples of ProGram being used to explain a demonstration grammar given in the appendix. The examples will demonstrate head feature handling, metarules, foot features, and FCD's. The grammar includes control agreement and FCR's, but because of the simplicity

lexicon, these will not have any overt effect in the examples. The given is more or less correct, in that it does what the author intended and in the examples ProGram will be used to explore different aspects rather than attempt to debug it further. There are, of course, inadequacies and gaps in its coverage, even of the small subset permitted by its lexicon.

First of all, a brief description of the grammar. The specifications follow GP82 fairly closely. There are four MAJOR feature categories, corresponding to verbal (V), nominal (N), prepositional (P), and determiner (D), and the grammar employs a simple three level X-bar system. Head features are AGR, which is used for the CAP, but will not be used here (all lexical items are singular), and VFORM, PFORM, NFORM, which are mutually exclusive, and specify type-specific data - for verbs, active, passive, for prepositions, the actual preposition (like a terminal feature) and for nouns, nothing. ~NFORM is used for determining simplicity. Three FCR's enforce the correct MAJOR/MINOR feature correspondence. There are nine ID rules and three metarules as described above, and in expansion, the grammar contains twenty rules derived as follows:

<u>Initial rules</u>	<u>One metarule</u>	<u>Two metarules</u>
nom	relcl(nom)	stm1(pass(vp2))
np2	stm1(pp)	pass(stm1(vp3))
np1	pass(vp3)	stm1(pass(vp3))
pp	stm1(vp3)	stm1(pass(vp3))
vp4	stm1(vp3)	
vp3	pass(vp2)	
vp2	stm1(vp2)	
vp1		
s		

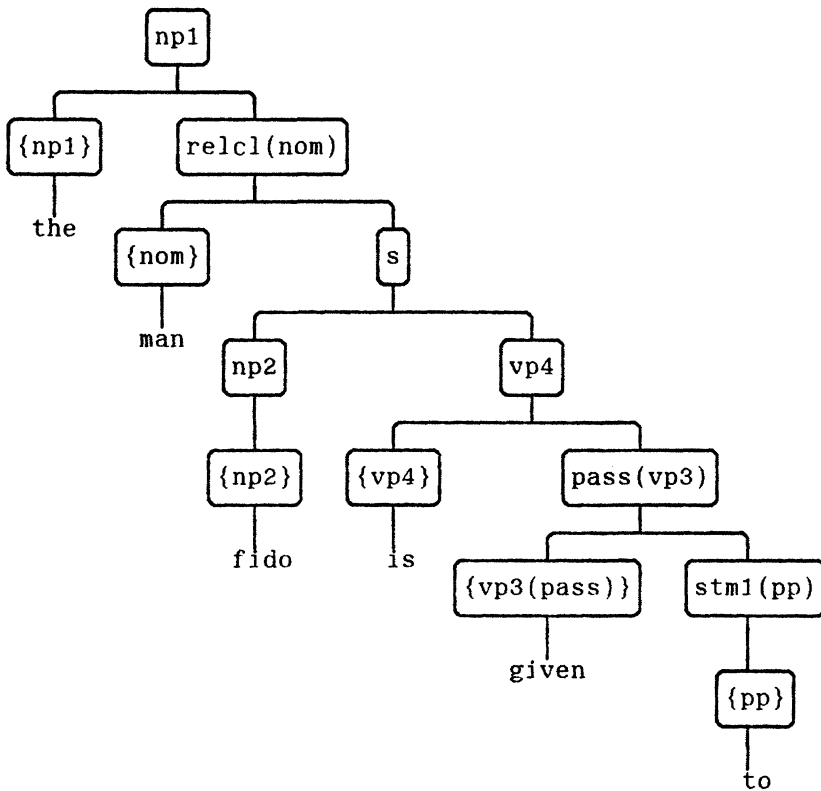
Notice that STM1 applies in two ways to VP3 (slashing either N(2) or P(2,to) and PASS(VP3) (slashing either P(2,to) or the optional N(2)). Unfortunately, this gives rules with the same names, which may cause confusion. The linear precedence rules are quite restrictive (English is strictly ordered): the only scope for variation is in VP3 where the P(2,to) are unordered. Thus the grammar accepts both "give fido to kim" and "give to kim fido". Finally the single FCD blocks passive verbs from being introduced except where they are specifically requested, namely in rules output from PASS.

Example 1: the FCD.

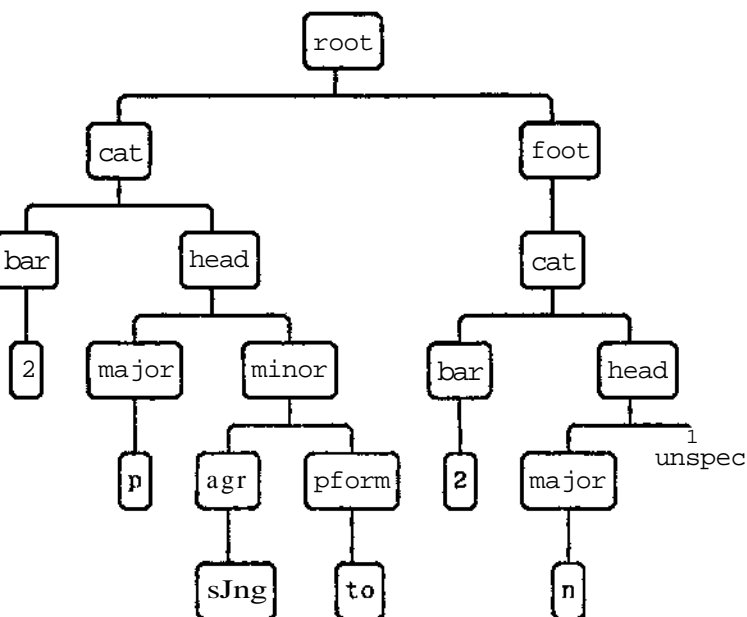
The single word SEEN has three parses, as a simple verb, as a passive verb phrase (as in "kim is seen") and as a passive verb phrase with a prepositional phrase in "by fido kim is seen", although the grammar does not handle topic markers. The FCD requires that SEEN cannot occur in the active verb position in automatic mode, a phrase such as "kim seen fido" has no parses. In control mode, ProGram announces that it is the FCD which ensures this.

Parsing in CONTROL mode.

```
|: kim seen fido                                The initial sentence to parse
--- LEX rule for FIDO |: np2                        Select the lexical rule for FIDO
    Located LEX rule np2
```

There is a departure from the usual conventions here: nodes are labeled by the mother category, but by the name of the ID rule applied. Thus, the categories are themselves feature trees and so cannot easily be represented in a compact form. Careful labelling of ID rules, however, ensures that the labels are sensible. Notice in the tree above that instances of only three metarules occur, STM1(PP), PASS(VP3) and RELCL(NOM). It is obvious from this tree is that four of the nodes have a slashed NP label: the S by the relative clause rule and consumed by the STM1(PP) rule, one needs to look at the feature structure of a category. ProGram can return the feature tree of any node on request. As an example, the node associated with a PP/NP category whose precise structure is as follows:



Notice that the MINOR features of the slashed NP were never specified, denoted in the tree by UNSPEC. It is also clear that attempting to use such a tree as a code label is rather impractical.

Example 3 - the Foot Feature Principle.

The parse tree given above is the only parse for the phrase "the man fido is given to" according to the example grammar. However, its sentential counterpart, "fido is given to the man", has two parses, as a straightforward sentence, and as a sentence with a slashed PP (the optional PP[by] is >ASS(VP3)). This second parse does not occur in the nominal form because the foot feature principle blocks it. It requires that "given" be consumed by JTM1(PASS(VP3)) as a VP/PP[by]. This rule does not require an object NP (PASS is deleted that) but it does need a PP[to]. In this phrase, there is a PP[to] but it has a slashed NP as well. The foot feature principle must transfer this NP to the mother, but it is incompatible with the /PP[by] already there. Thus the principle is violated. The following trace starts after the PP/NP ("to" has been built, and "given" is being introduced.

```

--- LEX rule for GIVEN | : VF3
    Located LEX rule vp3(pass)
  
```

```

Select VP3 for GIVEN
  
```

Current segment: (vp3(pass)} stml(pp)

——ID rule to consume {vp3(pass)}> |: vp3 Select a VP3 rule to build a VP
 Located ID rule vp3
- Ok? |: n There are six VP3 rules, we fore
 Located ID rule pass vp3 the parser to try a particular
- Ok? |: n one (in fact, the last one).
 Located ID rule stml vp3
- Ok? |: n (reject all these)
 Located ID rule stml vp3
- Ok? |: n
 Located ID rule pass stml vp3
- Ok? |: n
 Located ID rule stml pass vp3 This is the one.
- Ok? |: y

Warning: FFP failed.

Involving: [foot, [cat, [bar, [2]], [head, [major, [p]], [minor, J, [pform,
 [by]]]]]]]

*** No (more) possible matches for rule.
*** ID rule stml pass vp3 not applicable.

6> How ProGram works

The aim of this section is to give a brief description of some aspects of the actual mechanics of the ProGram system. The bulk of ProGram is written in Prolog [Clocksin and Mellish 1981] and conforms to the DEC10 Prolog standard [Pereira et al. 1978]. Prolog was chosen as the implementation language because of its relative portability and popularity, particularly in Europe, and increasingly elsewhere [9]. The core DEC10 compatible system contains all the main facilities, but the POPLOG [Hardy 1982] version also includes extra Prolog libraries which redefine and augment some of these features, providing graphical displays and editing facilities which could not easily be added to the standard system. As in section 4, the discussion will occasionally refer to some features of Prolog, and brief explanations will be given where necessary, and the interested reader is referred to [Clocksin and Mellish 1981] for full details.

DATA NORMALISATION

As discussed in section 3, data normalisation transforms the grammar into a standardised internal format, translating aliases and carrying out superficial error checking, the examples show how valuable this extra layer of translation is. The central task of the normaliser is the basic translation of a category expression. The formal representation of a feature (as defined in GP82) is already a bit cumbersome for comfort. ProGram's internal representation is in fact slightly more cumbersome, because it requires that all coefficients of a feature be present, and in a fixed order. Omitted feature coefficients are represented by Prolog variables. Thus, given the feature specification of the examples above, the category in (6.1) (in its internal representation) has alternative forms as in (6.2), but a unique ProGram formal

n (6.3).

```
(6.1) [root [cat [bar 2] [head [major v]]]]
```

```
(6.2) [root [cat [head [major v]] [bar 2]]]
      [root [cat [bar [2]] [head [major [v]]]]]
```

```
(6.3) [root [cat [bar [2]] [head [major [v]] M]] F]
```

n (6.3) M and F are Prolog variables representing the absent MINOR and FOO features, and the ordering of coefficients is fixed to be the same as in the feature specification. Note also that features not explicitly defined in the feature specification are assumed to take no coefficients. Thus 2 and V have the implicit definitions in (6.4).

```
(6.4) feature [2].
      feature [v].
```

The reason for standardising features thus is to make comparing and matching features straightforward. Indeed, feature unification (as defined in GP82) becomes just Prolog unification. In order to achieve this standard form, ProGram picks apart the category structure given, interpreting the aliases appropriately, and rebuilds the category with the translated coefficients inserted in the correct order. In the process it checks for unknown and incorrectly positioned features etc..

Most of the other normalisation processes are re-arrangements of the data into a more convenient form. For example, the BAR coefficient is extracted from each category and stored separately, so that it is easy to tell whether the category is lexical. And in ID rules, the rule name is attached as an extra sub-feature value of LEXICAL (the lowest bar level). This is also done to the categories in the lexicon rules, so that lexical subcategorisation is otherwise automatic. A lexical category with LEXICAL feature VP1, say, will only match a rule category with the same LEXICAL value, that is, a rule category in an ID rule labelled VP1.

RULE PRE-PROCESSING

The pre-processor takes the original set of ID rules and produces a new set which have been closed under metarule expansion, and which include the constraints required by the head feature convention and the control agreement convention. It is this latter set that are used by the parser, which consequently can ignore these aspects of rule instantiation completely. The reason these particular aspects are dealt with in this way is a combination of efficiency and convenience.

Following Thompson [Thompson 1981], ProGram is designed on the view that pre-expansion of the ID rules using the metarules is more sensible than the fly expansion. The overheads of incrementally calculating and saving new ID rules are far outweighed by the simplicity of pre-expansion, and the attendant ability to examine the expansion process in isolation, before any parsing is carried out.

The internal form of the metarules makes ID rule expansion particularly easy: if an ID rule can be successfully matched against the pattern of a metarule, then Prolog variables in the pattern will be given values from the ID

specific code in the parser, makes the handling of these principles efficient and totally transparent to the parser itself.

PROGRAM'S PARSER

Of the various aspects of the GPSG formalism, there are three which particularly suggest that the context-free grammars generated by GPSG grammars are in general best suited to bottom-up parsing techniques. The basic issue is whether the mother category or the daughter categories are the best selector of the correct context-free rule to apply [11].

Unless the linear precedence rules impose a total ordering on daughters of every rule, there will be rules which permit some freedom among daughters. Such rules generate several context-free rules with identical mother categories, but different daughter-lists. Thus a GPSG grammar which makes full use of ID/LP descriptions will contain in its underlying context-free grammar many groups of rules with the same mother category, and so mother category will be a bad selector of the correct context-free rule. Conversely, however, ID/LP does not significantly affect the extent to which a given daughter set selects a context-free rule.

Similarly, the Foot Feature Principle in GPSG requires that the mother's foot increment is the unification of the daughters' foot increments. In the details of this statement are not important here, the basic point to notice is that part of the mother is specified by unifying (an operation similar to intersection) parts of the daughters. Like ID/LP expansion, unification is an operation where many different configurations of 'inputs' (from daughters) specify the same 'output' (on the mother), so the problem described in the preceding paragraph arises again here.

Finally, the Conjunct Realisation Principle requires that the daughters of a coordinate structure are all extensions of the mother, that is, they are just more detailed instances of the same category as the mother. Here the mother could be any category contained in (that is, less well specified than) the maximal intersection of a given daughter set, and this means that in general, the daughters do not fully define the mother. However, ProG interprets the CRP as also requiring the mother to be maximal, and this means that the same considerations applied above are relevant here - the daughters specify the mother but not vice versa.

These conclusions suggest that bottom-up parsing is generally the best way to proceed with a GPSG grammar (whether fully expanded or not). The fact that the expansion (to a context-free grammar) has not been carried out in ProG provides further support. The categories in the rules are only fully specified once they have been incorporated into an analysis. A top-down approach must deal with this under-specification by guesswork, or by delayed instantiation (and hence delaying validity checks etc.) until after the daughter categories have been built. In either case it is difficult to avoid wasting time by following paths that cannot succeed due to earlier but as yet unchecked violations. An alternative which overcomes this problem is to carry out the checking incrementally, an approach which would require sophisticated inference mechanisms of the sort described in [Frisch 1985], but which are readily available in current programming languages. A bottom-up parser, on the other hand, starts off with fully instantiated lexical categories, and with care can ensure that categories are always fully specified and checked before they are used elsewhere.

Taking account of these arguments, ProGram parses bottom-up. The is very primitive, with no memory of partial constituents etc. It with the fully instantiated lexical categories, and attempts to find which will consume some of them. It takes an ID rule as its basis, and it against some of the daughters, checking linear precedence, and conventions. After a successful match, the mother category itself complete and so, inductively, all daughters are always complete. particular, it is guaranteed that the CRP and FFP will never modify in any way and so the linear precedence checks can be carried out as daughter is consumed by the ID rule, rather than waiting until the been found and FFP and CRP checks (which notionally precede the LP checks) been done. This rule matching process is then repeated, the overall being a single rule which consumes all the current categories.

The fact that the simplicity of the algorithm permits such optimal linear precedence etc. suggests that a more sophisticated algorithm have to be considerably more complex. The main reasons for keeping simple if inefficient schemes are that it is easy to follow when constructing a parser manually, and, because it uses Prolog's backtracking to do searching, it minimises its space requirement. The structures manipulated can be quite large, so this is a real consideration.

To give a brief idea of the algorithm, consider a sequence of as in (6.11) being parsed according to the simple grammar in (6.12). and subsequently, the bar represents the left boundary of the currently under consideration. The parser starts by considering the category.

(6.11) det n vtr det adj | n

(6.12) s --> np, vp
 np --> det, nom
 nom --> n
 nom --> adj, nom
 vp --> v, np

The parser will apply a rule to the N producing a NOM, and fail to apply to that (alone). So instead the parser moves the bar left to look at category (6.13).

(6.13) det n vtr det | adj nom

The rule consuming the ADJ and the NOM is now applied, creating a again, the NOM cannot be built up any further, so another category is the current list (6.14).

(6.14) det n vtr | det nom

From here, an NP can be built (but nothing more), the bar moved left built. The bar moves again, giving (6.15).

(6.15) det | n vp

This time, a rule applies to the N which does not consume the VP, so remains untouched. Again, no further building is possible, the bar and the NP and then the S can be built. With only one category remaining parse has been found.

At each stage of this procedure, there is one principle choice: whether to apply a rule (and if so which one) or move the bar left, to include another category. It is the parser's behaviour at such choice points which differs for different control modes. Exhaustive search in AUTO mode is achieved by using Prolog's backtracking - Prolog remembers each choice selected, and can back up to the last choice point and try a different alternative. Once all the rules and moving the bar left, have been attempted, Prolog backs up to the choice before and tries a different alternative there. In MONITOR mode, the user filters the choices the parser makes, and in CONTROL mode, the choice is left entirely up to the user. In the example, after producing the parse, the system does the last choice (building the S) and tries something else. The situation is as in (6.16).

(6.16) | np vp

There are no more rules to apply, and the bar cannot go further left, so the system backs up again, undoing the NP giving (6.17). Again nothing can be done and the next back up moves the bar RIGHT to give (6.18).

(6.17) | det nom vp

(6.18) det | nom vp

In this way the parser will back up all the way back to the beginning. Since the grammar is so simple, no new alternatives will be tried. In general, there may be new possibilities to try at any stage, causing the parser to move forwards again.

RULE MATCHING

The principle task of the parsing algorithm described above is the rule matching, which may be summarised as follows. The parser is faced with a list of categories (all those to the right of the bar) and a candidate ID rule. The daughter categories of the ID rule are to be matched with a leading (leftmost) sublist of the given categories such that

- a) linear precedence constraints are observed
- b) FCD specifications are fulfilled
- c) CRP and FFP between the mother and daughters are satisfied
- d) No essential daughters in the rule are not matched

If all these conditions hold, the rule has matched successfully, and the mother can replace the daughters in the list of given categories. Recall from above that HFC, CAP and metarule processing can be ignored.

An ID rule contains a list of daughter categories which are unordered and may be annotated with operators "opt" (optional category), "*" (zero or more instances of this category permitted) or "+" (one or more instances of this category permitted). Multiple instances need not be adjacent in the final rule (unless LP rules force them to) and if a multiple instance category is not fully instantiated, then each instance can be matched independently of the others. Such a rule is exhausted by a list of instantiated categories if every essential category in the rule matches one in the list (and possibly some non-essential ones do too). The basic ID rule matcher takes an ID rule and tries to allocate instantiated daughters (taken from left to right) to rule daughters (in any order) until there are no essential rule daughters outstanding. As each daughter is matched LP checks between it and all the daughters already matched

are carried out (it must not be required to precede any of them), sufficient to ensure LP admissibility of the entire list. Rule daughter given a status marker, which is suitably updated (according to their type) as they are consumed and used to check for rule exhaustion statuses are "requiring matching", "available for matching, but not matching" and "not available for matching. A rule is exhausted when there are no categories still requiring matching.

When the rule is exhausted, the CRP and FFP processing is carried out, complete the instantiation of the mother, and then FCR constraints are checked. FCD checks on the mother are deferred: since they depend on the category as daughter in a rule, they cannot be carried out until they are incorporated into the level above in the analysis.

FEATURE INSTANTIATION PRINCIPLES

The actual implementation of the remaining principles (FFP, CRP, and RAC) is straightforward. For the FFP, the daughter foot is collected up as the rule-match proceeds, and compared with the mother at the end. The CRP routines 'walk' down the feature trees of all the daughters in parallel, and as long as they coincide, the mother must agree with them. Checking takes place at the same time as rule category matching. Categories to be matched are examined in parallel to ensure that they agree with, but does not contribute to, the daughter, and non-matching features are checked for defaults. FCR and RAC checks are simple detection of particular feature value combinations on the now fully instantiated daughter.

7. Evaluation and Conclusions

This section provides a brief appraisal of ProGram, eight months after completion. Although the system has found its way into several sites in several countries in that time, there has as yet been little feedback from those trying to use it. Thus the comments in this section are based largely on the author's own experience of using and demonstrating the system, and the experiences of students and colleagues at Sussex [12].

The idea of using a controllable parser to 'debug' a grammar was provided as a tool which can show up problems with grammars that are not located manually. The overall design of the system sacrifices much of the interests of practicality. The segregation of data into different files, the disc-based nature of the system, and indeed the whole normalisation process could be more streamlined, or even completely transparent. To a certain extent ProGram could be improved as it is - the addition of more normalisation commands, and other high-level commands, would undoubtedly make the system easier to use, although probably no faster. But the overall automatic normalisation and sufficiently flexible file management would be a greater.

In general, the attempts to overcome the time and space problems have been partially successful. ProGram is a large system and it takes a great amount of skill to avoid loading too much of the system at once, but grammar development proceeds smoothly and quickly, and without inconveniencing other users or consuming processing quotas. However, with practice the system can be acquired, after which the system is in general quite manageable. Automatic parsing is a particular problem, and cannot always be avoided. The only remedies seem to be: keep phrases as short as possible (for efficiency).

lot try to test large noun phrases inside complex sentences if it can be avoided - use proper nouns instead), make sure no words are duplicated in the lexicon, keep the number of ID rules currently loaded down where possible, and run the parser in batch mode - perhaps using the cut-down parser, LIB PRS rather than the full system.

Understanding the behaviour of the parser in MONITOR and CONTROL mode also requires some effort. Most of the problem seems to be the right-to-left behaviour (which can be picked up quite quickly), and the backtracking/backtracking can produce some rather strange effects at times, although ProGram's behaviour is model Prolog chronological backtracking (no cuts etc.) In control mode, ProGram announces when it is forced to backtrack, which makes the flow of processing easier to follow.

ProGram provides a variety of tracing and diagnostic facilities which can all be very useful in various situations. There are, however, serious Missions, particularly in the CONTROL mode diagnostics. For example, it is not possible to look at the detailed structure of the categories currently under consideration during the parse. The WATCH flag provides parse tree information but this seems far less useful. One of the least helpful diagnostics is the one indicating simply that the rule match failed ('ID rule not applicable'). This can happen because the rule was totally inappropriate in the first place, because the output to some metarule was slightly incorrect, or because HFC or CAP have been violated. Because these latter features of GPSG are handled in the preprocessing stage, the parser simply cannot tell which is the problem. But if the categories themselves were available, it would often be possible to see what was wrong. More generally it is sometimes difficult for naive users to work out what the error messages mean: a given linear precedence problem, for example, might be caused by incorrect LP, ID rule or alias data (to mention but three), and the parser does not help decide which.

The lexicon seems to cause users some trouble. ProGram expects only the most rudimentary of lexicons and places quite severe constraints on it: the lexical categories must be fully specified and it is inadvisable for a word to occur more than once. Used sensibly, with only one or two words for each lexical type, and different words for, say, different case nouns, it is manageable, although the phrases that can be produced may sometimes seem a little contrived. Unfortunately, users of the system find it difficult to think in these terms, and try to create large lexicons and omit features which can make any value (for example, CASE for proper nouns). This leads to disaster since ProGram forces such features to stay unspecified, so that the categories never match with any rule. A lexicon utility has been implemented by Keller [Keller 1984], but this has not yet been fully integrated into the ProGram framework. It makes specification of a more general lexicon easier, and is a useful aid to developing a lexicon suitable for use with ProGram, and subsequently extending it to provide more realistic coverage.

The problem with the lexicon has a parallel with the parser itself. The algorithm is rather sensitive to the structure of the ID rules. For example, it dislikes explicit optional categories, especially occurring as left-most laughter in a rule. It will try almost everything else before it gets round to attaching such a category. It is far better behaved if optionality is encoded in two separate rules, one containing the category, the other not. The system prefers grammars that basically resemble 'standard' GPSG grammars, and with these it is generally well-behaved.

The ProGram system was completed in April 1984, and since subject to only minor modification and extension. The project is active and so further development is not foreseen. ProGram was an attempt to provide a computer representation of a GPSG grammar. The project was successful, providing a tool which offers effective GPSG grammar writer. Aside from ProGram's value as a development tool, the existence of the parser ensures that the internal formalism adequately captures all the aspects of GPSG (as presented in the ProGram system implementation provides a proven base for other GPSG system implementations etc.).

- 1] ProGram is a suite of Prolog [Clocksin and Hellish 1981] programs developed using the POPLOG programming environment [Hardy 1982] at the University of Sussex, supported by the Social Science Research Council (UK), grant number HR 7829/1 to the University of Sussex (principal investigator: Gerald Gazdar). ProGram was designed and implemented by Roger Evans and Gerald Gazdar. ProGram is available free to interested academic users - see below for details.
- 2] GPSG also contains a formal semantic component, but this will not be discussed here - the ProGram system contains no semantics.
- 3] As with any developing theory, several different variants of GPSG coexist at one time. Below, we shall commit ourselves to one particular formulation.
- 4] Throughout this paper, list notation [...] will be used without commas separating items. This is tidy, but NOT in keeping with Prolog syntax. The grammar in the appendix, however, does conform to Prolog syntactical specifications.
- 5] We need not be concerned here about the possibility that the underlying grammar is not context-free (although a full listing of an infinite rule set would undoubtedly cause a problem). The more recent proposal [Gazdar et al. 1985] that a GPSG grammar does not specify a context-free grammar but rather it constrains the tree-set directly, is more interesting, and corresponds more closely to the way ProGram actually works. However this difference is not important to the present discussion.
- 6] These conventions, and others mentioned below, are part of the GPSG theory and will be used without comment throughout this paper. The following abbreviations will be used: HFC - head feature convention, CAP - control agreement principle, FFP - foot feature principle, CRP - conjunct realisation principle (see [Gazdar et al. 1982]), DAC - default assignment convention, FCD - feature coefficient defaults (FSD in [Gazdar et al. 1985]), FCR - feature co-occurrence restrictions, RAC - root admissibility conditions. The reader is referred to GP82 for more details.
- 7] Aliases are in fact just ordinary Prolog clauses, so that the knowledgeable Prolog user can construct arbitrarily complex specifications. For example, the use of Prolog variables is particularly useful. The aliases in (4.14) are both subsumed by a single alias:

```
alias( v(N), [root [cat [bar N] [head [major v]]]]).
```

Here N is a Prolog variable (adopting the convention that Prolog variables start with an UPPER CASE letter), and the alias is valid for any value of N, but the same value must occur in both instances of N. (14) contain examples of this general alias with N=1 and N=2 respectively. See the example grammar for many further instances of this mechanism.
- 8] The following examples are real transcripts of ProGram use, slightly modified to remove unnecessary output etc. Input from the user is underlined.
- 9] ProGram has been tested with varying degrees of thoroughness on DEC1

Prolog, CProlog and POPLOG Prolog systems.

- [10] The Prolog notation [np|W] means "a list whose first element is np whose other (currently unspecified) elements are in the variable W, which is matched with a list like [np,v,s], W would be set to [v,s]."
- [11] This distinction is somewhat simplistic: for example, the properties of the first daughter in a context-free rule will be different in a predictive bottom-up parser. The general argument still holds, however, and the underlying question of what specific parsing algorithms are well suited to GPSG grammars (and hence, one might claim, to natural languages), has yet to be answered.
- [12] The author is indebted to the experience and comments of ProGram users at the University of Sussex, particularly Lynne Cahill and David Allport.

Availability

1. "The ProGram Manual" [Evans & Gazdar 1984] provides an introduction to the system, its features, and its use. It is available from the University of Sussex Cognitive Science Research Paper 35 (CSRP 35). It can be ordered from Ms. Judith Dennison, Cognitive Studies Programme, University of Sussex, Falmer, Brighton BN1 9QN, for 7.50 pounds including postage and packing.
2. ProGram is now part of the standard Sussex POPLOG system and is available without extra charge, in all academic issues and updates of the POPLOG system. POPLOG is available for VAX's under VMS, VAX's under MS-DOS, and Bleasdale BDC 680a's under Unix. Non-educational customers (UK & elsewhere) who want ProGram with POPLOG should order it through System Software Ltd., Systems House, 1 Pembroke Broadway, Camberley, Surrey GU15 2AA (0443 62244).
3. Academic users of other Prolog systems can obtain a magnetic tape (in tar format) of the Prolog code of the ProGram system free, together with a copy of "The ProGram Manual", provided they pay the tape, postage, and handling costs (35 pounds). Copies can be ordered from Alison Mudd, Cognitive Studies Programme, Arts E, University of Brighton, Brighton BN1 9QN. A cheque for 35 pounds, made payable to "The University of Sussex", should be enclosed with the order.

References

- Clocksini and Mellish 1981: Programming in Prolog. Clocksin W. and C. Mellish. Springer Verlag, 1981.
- Evans and Gazdar 1984: The ProGram Manual. Evans R. and G. Gazdar. Cognitive Studies Research Paper No: CSRP 035, University of Sussex, April 1984.
- Farkas et al. 1983: Some revisions to the theory of features and instantiation. Farkas D., D. Flickinger, G. Gazdar, W. Ladusaw, J. Pinkham, G. Pullum and P. Sells. In Proceedings of the ICOT Workshop on Non-Transformational Grammars, 11-13 (Tokyo: Institute for New Computer Technology). 1983

- Frisch 1985: Parsing with Restricted Quantification. A. Frisch. University of Sussex, January 1985
- Gazdar et al. 1982: Coordinate structure and unbounded dependencies. Gazdar G., E. Klein, G. Pullum, and I. Sag. In M. Barlow, D. Flickinger & I.A. Sag (eds.) Developments in Generalised Phrase Structure Grammar: Stanford Working Papers in Grammatical Theory, Volume 2. Bloomington: Indiana University Linguistics Club, 38-68. 1982
- Gazdar et al. 1985: Generalised Phrase Structure Grammar. Gazdar G., E. Klein, G. Pullum, and I. Sag Oxford: Blackwell, 1985.
- Gazdar and Pullum 1982 (GP82): Generalised phrase structure grammar: theoretical synopsis. Gazdar G. and G. Pullum. Mimeo, Indiana University Linguistics Club, August 1982.
- Hardy 1982: The POPLOG Programming System. Hardy S., Cognitive Studies Research Paper No: CSR P 003, University of Sussex, November 1982
- Keller 1984: A Lexicon Handler for the ProGram Grammar Development System. W. R. Keller Cognitive Studies Research Paper No: CSR P 040, University of Sussex, June 1984
- Pereira et al. 1978: User's Guide to DECsystem-10 prolog. Pereira L. M., F. C. N. Pereira and D.H.D Warren. DAI Occasional paper 15, Department of Artificial Intelligence, University of Edinburgh. 1978
- Pereira and Warren 1980: Definite Clause Grammars for Language Analysis - a survey of the formalism and a comparison with augmented transition networks. Pereira F. C. N., and D.H.D Warren. Artificial Intelligence 13:3 231-278. 1980
- Thompson 1981: Handling metarules in a parser for GPSG. Edinburgh D.A.I. Research Paper No. 175. Also: In M. Barlow, D. Flickinger & I.A. Sag (eds.) Developments in Generalised Phrase Structure Grammar: Stanford Working Papers in Grammatical Theory, Volume 2. Bloomington: Indiana University Linguistics Club, 26-37.

Appendix - the example grammar

This appendix contains the complete grammar used in the paper. Note version below conforms to Prolog syntax, (unlike the examples in text).

Features

syncat root.

```
feature [root, cat, foot].
feature [cat, bar, head].
feature [bar, {lexical, 1, 2}].
feature [head, major, minor].
feature [major, {v, n, d, p>}].
feature [minor, agr, {vform, pform, nform}].
feature [agr, {sing, plur}].
feature [vform, {fin, pass}].
feature [pform, {by, to}].
feature [foot, cat].
```

Aliases

```
alias( v(N), [root,[cat,[bar,N],[head,[major,v]]]] ).
alias( n(N), [root,[cat,[bar,N],[head,[major,n]]]] ).
alias( d(N), [root,[cat,[bar,N],[head,[major,d]]]] ).
alias( p(N), [root,[cat,[bar,N],[head,[major,p]]]] ).
alias( h(N), [root,[cat,[bar,N],[head,[major]]]] ).
alias( v(L,N),[root,[cat,[bar,L],[head,[major,v],[minor,[vform,N]]]] ).
alias( p(L,N),[root,[cat,[bar,L],[head,[major,p],[minor,[pform,N]]]] ).
alias( h, h(lexical) ).
alias(X/Y, Z) :-
    normfeat(X,XN),normfeat(Y,YN),
    pathfor(foot,YN,f~'),
    pathfor(cat,YN,YCat),
    pathfor(foot,XN,[[cat]YCat]),
    Z = protect(XN).
```

(The following aliases are mainly for the lexicon:)

```
alias( nO, [root,[cat,[bar,lexical],
                [head,[major,n], [minor,[agr,sing],nform]]] ),
alias( dO, [root,[cat,[bar,lexical],
                [head,[major,d], [minor,[agr,sing], "nform]]] ),
alias( vO(V), [root,[cat,[bar,lexical],
                [head,[major,v], [minor,[agr,sing],[vform,
alias( pO(P), [root,[cat,[bar,lexical],
                [head,[major,p], [minor,[agr,sing],[pform,
```

ID rules

```
s: v(2) -> N2,H1 where N2 is n(2),
                        HI is h(1),
                        N2 controls HI.
```


nom:	n0	->-	woman, man.
np1:	d0	->-	the.
pp:	p0(to)	->-	to.
pp:	p0(by)	->-	by.