

THE COMPUTER STORAGE, RETRIEVAL AND SEARCHING

OF GENERIC STRUCTURES IN CHEMICAL PATENTS:

THE MACHINE-READABLE REPRESENTATION

OF GENERIC STRUCTURES

A Study Submitted in Fulfilment of the Requirements

for the Degree of

DOCTOR OF PHILOSOPHY

by

JOHN MORDAUNT BARNARD

December 1982

Department of Information Studies

University of Sheffield

## CONTENTS

CHAPTER 1 : GENERIC STRUCTURES IN PATENTS.....	1
1.1. THE NATURE OF GENERIC STRUCTURES.....	3
1.1.1. Patent Claims.....	3
1.1.2. Types of Generic Structure.....	5
1.1.3. Generic Structures Outside Patents.....	7
1.2. GENERIC STRUCTURE DESCRIPTIONS IN PATENTS.....	8
1.2.1. The Constant Part.....	8
1.2.2. The Variable Parts.....	9
1.2.3. Conditional Expressions.....	13
1.3. THE "MARKUSH PROBLEM".....	13
1.4. GENERIC STRUCTURE REPRESENTATIONS.....	15
1.4.1. Derwent Publications Ltd.....	16
1.4.2. IFI/Plenum Data Co.....	19
1.4.3. International Documentation in Chemistry.....	21
1.4.4. Chemical Abstracts Service.....	25
1.4.5. Système DARC.....	27
1.4.6. Line Notations.....	28
1.4.7. The COUSIN system.....	32
1.5. REQUIREMENTS FOR A SEARCH SYSTEM REPRESENTATION.....	33
CHAPTER 2 : FORMAL LANGUAGES.....	37
2.1. DEFINITION AND CLASSIFICATION OF FORMAL LANGUAGES.....	38
2.1.1. The Chomsky Hierarchy.....	41
2.2. PARSING OF CONTEXT-FREE LANGUAGES.....	43
2.2.1. LR Parsing.....	45
2.2.2. LL Parsing.....	46
2.2.3. Top-Down vs. Bottom-Up Parsing.....	47
2.3. PROGRAMMING LANGUAGES.....	48
2.3.1. Syntax Specification.....	49
2.3.2. Syntactic Analysis.....	51
2.3.3. The Pascal Language.....	52
2.3.4. The Ada Language.....	54
2.3.5. Choice of Language for Software Development.....	55
2.4. FORMAL LANGUAGE SEMANTICS.....	56

## CONTENTS

2.5. INTERACTIVE LANGUAGES.....	58
2.6. FORMAL LANGUAGES IN CHEMISTRY AND INFORMATION WORK.....	60
CHAPTER 3 : THE INPUT LANGUAGE.....	64
3.1. GENERIC STRUCTURE DESCRIPTION USING GENSAL.....	65
3.2. STRUCTURE DIAGRAM INPUT.....	68
3.3. SIMPLE ASSIGNMENT STATEMENTS.....	71
3.3.1. Substituent Assignments.....	72
3.3.1.1 Unknown Value.....	73
3.3.1.2 Structure Diagram.....	73
3.3.1.3 Nomenclatural Terms and Expressions.....	73
3.3.2. Multiplier Assignments.....	75
3.4. MORE COMPLEX ASSIGNMENTS.....	76
3.4.1. Combined Substituents.....	77
3.4.2. Group Assignment Statements.....	78
3.4.2.1 Assignment Operators.....	79
3.4.2.2 Selected Group Assignments.....	82
3.5. HOMOLOGOUS SERIES IDENTIFIERS AND GRAMMARS.....	82
3.6. POSITION SETS.....	86
3.7. NESTED SUBSTITUTION.....	89
3.7.1. Selectors in Definition Expressions.....	91
3.7.2. Position Sets in Definition Expressions.....	92
3.7.3. Substituents as Substituent Values.....	92
3.7.4. Further Substitution on Parenthesised Expressions.....	94
3.8. SPECIAL RESTRICTIONS IN GENERIC STRUCTURES.....	96
3.8.1. Conditions.....	97
3.8.2. Definition Relations.....	98
3.8.3. Integer Relations.....	99
3.8.4. Group Relations.....	101
3.8.5. IF Statements.....	102
3.8.6. RESTRICT Statements.....	104
3.9. SCOPE OF DEFINITIONS.....	105
3.10. LIMITATIONS OF GENSAL.....	105
3.11. THE DESIGN OF GENSAL.....	108
3.11.1. Formal Grammar.....	108
3.11.2. Non-Determinacy.....	109
3.11.3. Security vs. Flexibility.....	111

## CONTENTS

CHAPTER 4 : THE INTERNAL REPRESENTATION.....	114
4.1. REQUIREMENTS FOR THE REPRESENTATION.....	115
4.2. THE PARTIAL STRUCTURE RECORD.....	119
4.2.1. Specific Partial Structures.....	120
4.2.2. Generic Partial Structures.....	120
4.2.3. Unknown Partial Structures.....	121
4.2.4. Other Partial Structures.....	121
4.3. CONNECTION TABLE FORMAT.....	122
4.3.1. Congener Record.....	123
4.3.2. Bond Orders.....	124
4.4. PARAMETER LIST FORMAT.....	125
4.5. CHILD GATE FORMAT.....	126
4.5.1. Combination Bars.....	128
4.5.2. Alternative Bars.....	130
4.6. PARENT GATE FORMAT.....	130
4.7. REPRESENTATION OF CONDITIONS AND RESTRICTIONS.....	132
4.8. THE ECTR AND OTHER REPRESENTATIONS.....	134
4.9. IMPLEMENTATION OF THE ECTR.....	136
4.9.1. The Partial Structure Record.....	136
4.9.1.1 Connection Tables.....	137
4.9.1.2 Parameter Lists.....	139
4.9.1.3 Other Terms.....	140
4.9.2. Child Gate Record.....	140
4.9.3. Parent Gate Record.....	142
4.9.4. Space Requirements.....	143
CHAPTER 5 : AN INTERPRETER FOR GENSAL.....	145
5.1. INVOCATION OF THE INTERPRETER.....	147
5.2. LEXICAL ANALYSIS.....	148
5.3. SYNTAX ANALYSIS.....	149
5.4. ERROR HANDLING.....	150
5.4.1. Program Errors.....	151

## CONTENTS

5.4.2. Structure Diagram Errors.....	151
5.4.3. "Immediate" Errors.....	151
5.4.4. "Delayed" Errors ("Failures").....	152
5.5. STRUCTURE DIAGRAM PROCESSING.....	152
5.5.1. The Feldmann Program.....	153
5.5.2. Procedure PROCESSCT.....	154
5.5.3. Storage of GENSAL Structure Diagrams.....	155
5.6. SUBSTITUENT DECLARATIONS.....	156
5.7. SUBSTITUENT DEFINITIONS AND ECTR GENERATION.....	158
5.7.1. Syntactic and Semantic Analysis in ELEMENT.....	159
5.7.2. Substituent Values.....	160
5.7.3. Nomenclatural Terms.....	161
5.7.4. Parameter Lists.....	165
5.7.5. ECTR Generation.....	166
5.8. MULTIPLIER DECLARATIONS AND DEFINITIONS.....	167
5.9. TIDYING THE ECTR.....	168
CHAPTER 6 : CONCLUSIONS AND SUGGESTIONS FOR FURTHER WORK.....	169
6.1. DEVELOPMENT OF A PATENT DOCUMENTATION SYSTEM.....	169
6.2. OTHER POTENTIAL APPLICATIONS OF GENSAL.....	172
6.2.1. Non-Computer Description of Generic Structures...172	
6.2.2. Generic Structures in the Journal Literature.....173	
6.2.3. Chemical Reaction Documentation.....173	
6.2.4. Specific Structure Search Queries.....174	
6.3. CONCLUSIONS.....	175
APPENDIX 1. GENSAL Syntax Diagrams.....	177
APPENDIX 2. BNF Syntax for GENSAL.....	182
APPENDIX 3. GENSAL Interpreter Program.....	186
APPENDIX 4. Program Index.....	335
APPENDIX 5. Global Declarations in GENPROG.....	339

**CONTENTS**

**APPENDIX 6. Sample Interpreter Session.....357**

**APPENDIX 7. Interpreter Error Messages.....363**

**BIBLIOGRAPHY.....367**

THE COMPUTER STORAGE, RETRIEVAL AND SEARCHING  
OF GENERIC STRUCTURES IN CHEMICAL PATENTS:  
THE MACHINE-READABLE REPRESENTATION OF GENERIC STRUCTURES

Thesis submitted for the Degree of Ph.D. by J.M. Barnard

**ABSTRACT**

The nature of the generic chemical structures found in patents is described, with a discussion of the types of statement commonly found in them. The available representations for such structures are reviewed, with particular note being given to the suitability of the representation for searching files of such structures. Requirements for the unambiguous representation of generic structures in an "ideal" storage and retrieval system are discussed.

The basic principles of the theory of formal languages are reviewed, with particular consideration being given to parsing methods for context-free languages. The Grammar and parsing of computer programming languages, as an example of artificial formal languages, is discussed. Applications of formal language theory to chemistry and information work are briefly reviewed.

GENSAL, a formal language for the unambiguous description of generic structures from patents, is presented. It is designed to be intelligible to a chemist or patent agent, yet sufficiently

## ABSTRACT

formalised to be amenable to computer analysis. Detailed description is given of the facilities it provides for generic structure representation, and there is discussion of its limitations and the principles behind its design.

A connection-table-based internal representation for generic structures, called an ECTR (Extended Connection Table Representation) is presented. It is designed to represent generic structures unambiguously, and to be generated automatically from structures encoded in GENSAL. It is compared to other proposed representations, and its implementation using data types of the programming language Pascal described.

An interpreter program which generates an ECTR from structures encoded in a subset of the GENSAL language is presented. The principles of its operation are described.

Possible applications of GENSAL outside the area of patent documentation are discussed, and suggestions made for further work on the development of a generic structure storage and retrieval system based on GENSAL and ECTRs.

## NOTE

The work described in this Thesis has been undertaken as part of a more comprehensive project on the computer storage and retrieval of generic chemical structures in patents. Whilst this has involved close liason with the other research worker on the project, S.M. Welford, the work described in this Thesis is entirely that of the author.

A number of publications have appeared describing work on the project; <sup>174-177</sup> the substance of Chapter 3 appeared in the second of these <sup>176</sup> and the substance of Chapter 4 in the third. 177

In addition, presentations have been given at the following meetings:

1. Chemical Notation Association (UK) Seminar on "Structure Searching in the Published Literature", Daresbury, March 1980.
2. Chemical Notation Association (UK) Seminar on "The Future of Chemical Documentation", Exeter, September 1982.

## ACKNOWLEDGEMENTS

I should like to thank my supervisor, Professor M.F. Lynch for his constant support and encouragement during the course of this research, and for inviting me to participate in it.

One of the satisfactions of this project has been the opportunities it has afforded to discuss my work with a wide range of individuals from different organisations, especially in the chemical and pharmaceutical industries. I should particularly like to thank the following for their many helpful comments and advice: John Silk (EUSIDIC, formerly ICI PPD), Peter Steele (Glaxo), Clive Tomlin, Richard Waterman, and David Pearson (all of ICI PPD), Frank Jackson (Pfizer), Charles Oppenheim and Peter Norton (Derwent), George Adamson (ICI Pharmaceuticals), Bill Town and Ole Norager (ISPRA), Todd Wipke (University of California), Claus Suhr (BASF) and Dr G.Pötscher (Fachinformatationszentrum Chemie, Berlin).

In addition, from my own Department I should like to thank Peter Willett, and Annette von Scholley who has still not managed to send my program down.

I should like to thank Mike Elder and Steve Hull (SERC, Daresbury) for providing the software for the Feldmann chemical structure graphics system, and for helping me to get it working.

## ACKNOWLEDGEMENTS

None of the programming work in this project would have been possible without the dedication to duty of the staff of the University Computing Services Department, and I should especially thank Richard Gilbert and Chris Martin for so efficiently removing all the bugs I found in their Pascal compiler. Also, from the Computer Science Department, I would like to thank Siobhan North for setting me along the right road in the matter of syntax analysis.

Finally, and most importantly, I should like to express my gratitude and appreciation to my friend and colleague Stephen Welford, who has spent countless hours discussing this project with me, and who also put up with sharing an office with me for two and a half years.

## CHAPTER 1

### GENERIC STRUCTURES IN PATENTS

---

"Bloody instructions, which, being taught  
return to plague the inventor"

Macbeth, Act I, Sc. vii

On the basis of those words, it might well be supposed that Macbeth was an information scientist in the chemical or pharmaceutical industries. Those industries are not only prolific generators of patent documents, but are also major users of patent information, and the efforts made to protect a company's invention in drafting a patent return to cause many problems for information searchers in the patent literature.

The increase in the number of chemical patent documents published in recent years has been prodigious, and the increase has continued in spite of a slight fall-off in the number of journal

articles published. In 1981 more than 71 000 patents were abstracted in Chemical Abstracts, as compared with fewer than 62 000 the previous year,<sup>1</sup> and this continues a trend which can be traced back many decades<sup>2</sup> though at least part of the increase can be explained by improvements in the range of countries covered by Chemical Abstracts; its relatively poor coverage compared to other indexing systems had previously attracted criticism.<sup>3</sup>

A further factor in the increase in published patent documents is the change in patent legislation in a number of countries, including Britain, during the 1970's.<sup>4, 5</sup> This has resulted in a move from the publication of examined and accepted patents only to the publication of unexamined applications. Initially this led to the sudden publication of backlogs of applications, increasing the figures for patent documents published, but it has also led to a change in the actual substance of patent claims, especially in the chemical area, which has itself caused problems for patent documentation systems.<sup>6</sup>

This is because patents for chemicals and pharmaceuticals frequently do not lay claim to the single compound which the company taking out the patent intends to market, but rather lays claim to a whole class of compounds having broadly the same properties. In the initial application for a patent, a company may attempt to claim as wide a range of compounds as possible, partly to cover anything which might conceivably have the desired activity (patent application normally takes place well before

testing and development of "lead compounds" has been completed), partly to intimidate rival companies who may be working in the same area, and partly to disguise the true nature of the invention. It is very possible that the initial application may have to be modified before it can be accepted and a patent granted, but under the early publication system now adopted by most countries, it is the initial application which is published first. This retains its significance after examination - and many patent applications are in fact abandoned, no examination taking place and no patent being granted - as the information contained in it may affect the validity of future patents.

The class of compounds claimed in a patent is described by means of a generic structure which contains both fixed and variable parts, the extent of the variation defining the size of the class of structures.

## 1.1. THE NATURE OF GENERIC STRUCTURES

### 1.1.1. Patent Claims

In 1924 an American chemist, Eugene A. Markush, applied for a patent for a class of novel pyrazolone dyes,<sup>7</sup> but his application was rejected on the grounds that it claimed alternatives. After making suitable changes to the wording of his

application in order to leave out the word "or", it was accepted, and since then the term "Markush" has been applied to this type of generic structure. Rosa <sup>8</sup> has discussed the legal wrangles over this and other applications, and outlined the type of generic structure which may be claimed under the precedent set by Markush, though the rigid "Rule Against Or", <sup>9</sup> which never applied in other countries, has now been abolished in the United States too.

The expression "Markush Structure" is now used rather loosely to refer to a wide variety of types of generic structure, though U.S. patent attorneys use it to refer specifically to patents granted under the precedent established by Markush's pyrazolone dyes application. On account of this special legal meaning the expression has generally been avoided in the present work following advice from Silk <sup>10</sup> and despite its use by many other authors in the field, and the expression "generic structure" is used throughout this Thesis.

A single generic structure may cover an enormous, and in some cases infinite, number of specific compounds, <sup>11</sup> only a tiny fraction of which have actually been tested for the claimed activity. Beton <sup>12</sup> cites the example of a patent application on sulphathiazole which was rejected because, of the at least 93 million specific compounds covered, only two had been shown to have the claimed activity. In the same paper however, he refers to the original patent on the Ziegler process for ethylene polymerisation in which aluminium trialkyl is claimed as

catalyst. Following grant of this patent, Ziegler found that alkyl aluminium halides and organomagnesium compounds could also be used, and was obliged to make further applications to cover them also. However, this still left him with no patent protection for the use of such catalysts in the polymerisation of other alkenes.

These examples illustrate the need to formulate a patent specification sufficiently widely to cover all the compounds with the required activity, yet sufficiently narrowly not to claim untested compounds which are actually inactive.

### 1.1.2. Types of Generic Structure

Valance <sup>13</sup> has discussed the variety of different types of statement that may be found in generic structures, with a survey of their relative frequencies. Sneed, Turnipseed and Turpin <sup>14</sup> have attempted a rudimentary classification of generic structures, dividing them into **determinate** and **indeterminate** structures, the former having variable substructure groups (all defined) occurring with variable frequency at fully-defined positions of attachment, and the latter comprising all other generic structures, including those involving verbal expressions, undefined substructures and undefined positions of attachment. Concentrating on **determinate** structures, they give examples of the different types of expression that may be found.

A similar classification has been given by Krishnamurthy and Lynch <sup>15, 16</sup> dividing generic structures into delimited and undelimited structures, though these classes are not identical with Sneed et al.'s determinate and indeterminate structures. Delimited structures are essentially those which cover a finite (even if very large) number of specific compounds; undelimited those which cover an infinite number of specific compounds.

In the present work these classifications have not been found helpful, and analysis of generic structures has been based on an approach given by Geivandov <sup>17</sup> which views such a structure as a (possibly vestigial) constant part to which are attached variable parts that can vary in their chemical nature, in their position of attachment to the constant part, and in their multiplicity of occurrence. This concept may be extended to encompass the idea of a "Markush within a Markush" so that each variable part can have further variable parts attached to it, continuing to any level.

On this basis, two opposite "extremes" of generic structures may be identified: that where the "variable" parts are fully defined in terms of nature, position and multiplicity, in which case the structure is a specific structure identifying a unique chemical substance, and that where the variable parts are totally undefined, in which case the structure is a substructure which may be found embedded in any of a potentially infinite variety of specific structures.

Between these extremes lie generic structures with incompletely-

defined variable parts. Any variable part may still have an infinite number of different possible values, but it is none the less restricted in some way. For example, the term "alkyl" strictly covers the infinite variety of radicals containing carbon and hydrogen only, with no double or triple bonds and no rings, but it nonetheless restricts the variety of values a group defined as "alkyl" can take.

### 1.1.3. Generic Structures Outside Patents

Generic structures are also found outside patents. They appear in the journal literature, where a large number of related compounds have been tested for a particular property or activity, and in this case a generic structure is essentially a shorthand way of listing the compounds tested. Figure 1.1 shows an example of a generic structure from the Journal of Medicinal Chemistry.

Generic structures may be used as queries in some chemical structure search systems, with databases of specific structures. Generally, only very simple generic structures can be used, but the recently-developed COUSIN system allows more complicated queries. This is discussed more fully in Section 1.4.7.

The description of generalised chemical reactions can involve the use of generic structures for the generalised reactants and products, though no reaction indexing system has yet been developed using such reactant and product descriptions.

*Hallucinogenic Amphetamines* *Journal of Medicinal Chemistry, 1977, Vol. 20, No. 12 1633*

**Table I. Substituted Amphetamines and Predicted Hallucinogenic Activity**

No.	Ring position and group						<sup>3</sup> X <sub>D</sub>	<sup>4</sup> X <sub>D</sub>	<sup>5</sup> X <sub>pe</sub> <sup>v</sup>	Exptl <sup>a</sup> log μ	Calc <sup>d</sup> log μ
	2	3	4	5	6						
1			OCH <sub>3</sub>				3.348	1.034	0.469	0.59	0.55
2	OCH <sub>3</sub>		OCH <sub>3</sub>				4.124	1.508	0.642	0.67	0.87
3	OCH <sub>3</sub>			OCH <sub>3</sub>			4.124	1.683	0.638	0.87	1.06
4		OCH <sub>3</sub>	OCH <sub>3</sub>				4.808	1.830	0.739	0.37	0.55
5	OCH <sub>3</sub>	OCH <sub>3</sub>		OCH <sub>3</sub>			4.830	2.083	0.765	0.63	1.01
6	OCH <sub>3</sub>		OCH <sub>3</sub>		OCH <sub>3</sub>		4.853	2.031	0.798	1.03	1.12
7	OCH <sub>3</sub>	OCH <sub>3</sub>			OCH <sub>3</sub>		4.933	2.045	0.810	1.14	1.06
8	OCH <sub>3</sub>		OCH <sub>3</sub>	OCH <sub>3</sub>			4.892	2.058	0.785	1.26	1.00
9	OCH <sub>3</sub>	OCH <sub>3</sub>	OCH <sub>3</sub>	OCH <sub>3</sub>			5.629	2.363	0.917	0.86	0.92
10			-OCH <sub>2</sub> O-				4.203	1.705	0.576	0.41	0.21
11		OCH <sub>3</sub>		-OCH <sub>2</sub> O-			4.925	2.252	0.707	0.43	0.62
12		-OCH <sub>2</sub> O-	OCH <sub>3</sub>				5.043	2.272	0.756	0.48	0.80
13	OCH <sub>3</sub>		-OCH <sub>2</sub> O-				5.027	2.197	0.753	1.00	0.71
14	OCH <sub>3</sub>			-OCH <sub>2</sub> O-			4.993	2.317	0.751	1.08	0.71
15	OCH <sub>3</sub>	OCH <sub>3</sub>		-OCH <sub>2</sub> O-			5.746	2.749	0.685	0.75	1.09
16	OCH <sub>3</sub>		-OCH <sub>2</sub> O-		OCH <sub>3</sub>		5.761	2.906	0.887	1.13	1.29
17	OCH <sub>3</sub>			OC <sub>2</sub> H <sub>5</sub>	OCH <sub>3</sub>		5.027	2.285	0.768	1.22	0.93
18	OCH <sub>3</sub>			Br	OCH <sub>3</sub>		4.574	1.762	1.157	2.71	2.92
19	OCH <sub>3</sub>			CH <sub>3</sub>	OCH <sub>3</sub>		4.574	1.762	0.888	1.89	1.85
20	OCH <sub>3</sub>			C <sub>2</sub> H <sub>5</sub>	OCH <sub>3</sub>		4.892	2.058	0.910	2.01	1.70
21	OCH <sub>3</sub>			n-C <sub>4</sub> H <sub>9</sub>	OCH <sub>3</sub>		5.027	2.285	0.850	1.94	1.60
22	OCH <sub>3</sub>			n-C <sub>6</sub> H <sub>13</sub>	OCH <sub>3</sub>		5.296	2.436	0.880	1.63	1.34
23	OCH <sub>3</sub>			n-C <sub>8</sub> H <sub>17</sub>	OCH <sub>3</sub>		5.546	2.548	0.880	1.09	1.09

<sup>a</sup> Molar basis, ref 1 and 5, converted by multiplying by the ratio of molecular weights of amphetamine to mescaline.

Figure 1.1: A generic structure from the journal literature

## 1.2. GENERIC STRUCTURE DESCRIPTIONS IN PATENTS

The manner of description of generic structures in patents from different countries is basically equivalent, and an example of such a description from a recent British patent is shown in Figure 1.2.

### 1.2.1. The Constant Part

In a typical patent specification, or abstract, there is a structure diagram for the constant part, in which the attached variable parts are indicated by symbols such as R, X, R'', R<sub>2</sub> etc. There is little or no standardisation of the symbols used, and

occasionally valid atomic symbols (such as B or C) appear as structural variables, which can cause ambiguity and confusion.

The variables may be attached to the constant part at fixed or variable positions, the latter normally being indicated by the convention of a bond going into the centre of a ring, or sometimes, where the attachment is to a chain, by means of a brace over the possible atoms of attachment. The variables may have one or two connections to the constant part, with any bond orders, or infrequently three or more.

Multiplicity of occurrence of certain portions (normally structural variables) of the structure diagram is often indicated by a subscript to a symbol, or to parentheses around a multiplied portion. The subscript may be a single integer, a range of integers, or an alphabetical or other symbol that is defined elsewhere. Examples are:



### 1.2.2. The Variable Parts

Following the constant part, the variables introduced in it are defined, usually by listing the alternative values for each structural variable. However, several different types of

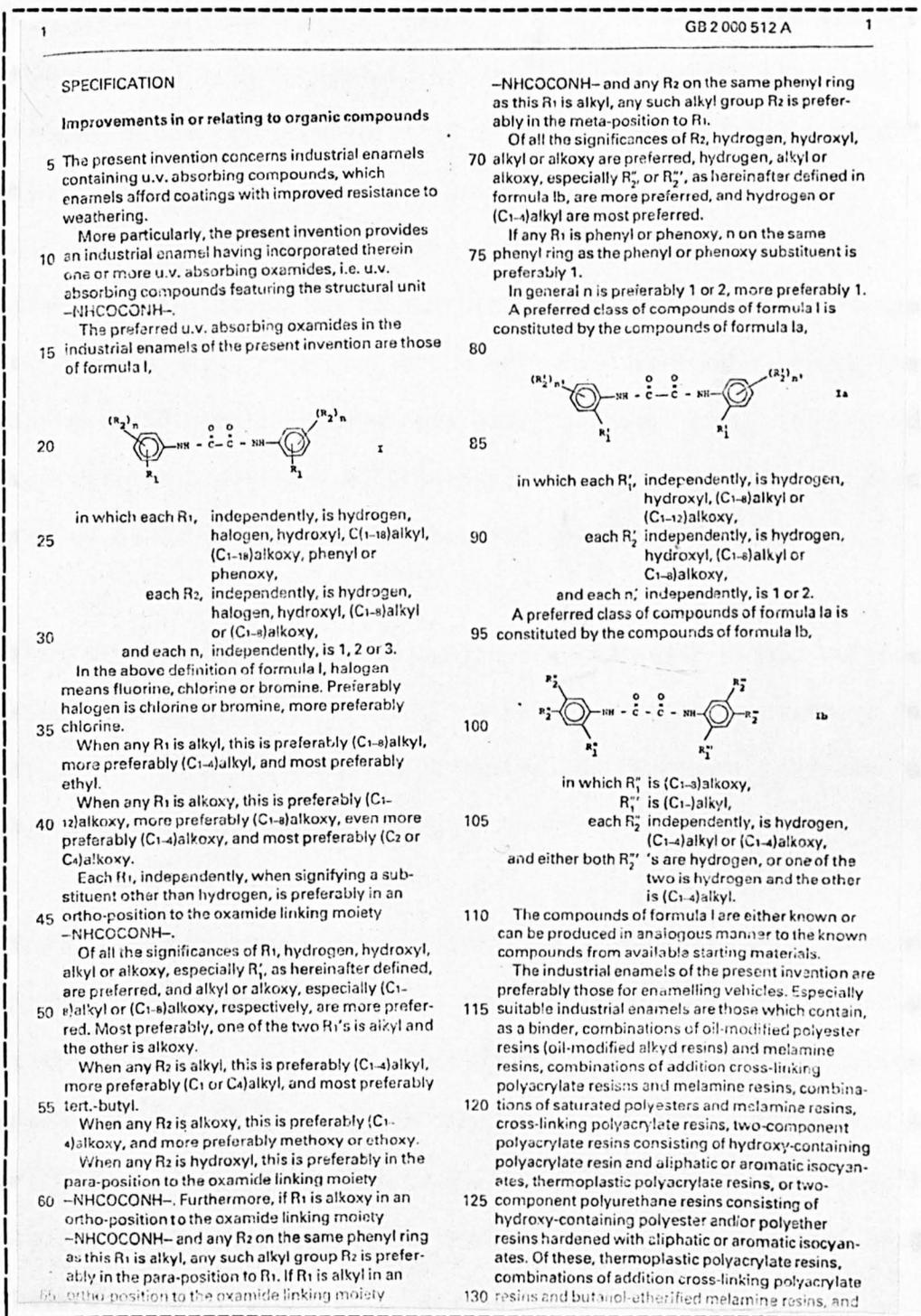


Figure 1.2: Part of a British Patent Specification

expression may be used for these values.

There may be simple nomenclatural terms (e.g. "methyl",

"cyclohexyl", "pyridyl", "amino", etc.) that represent single chemical entities, or terms or expressions that represent a limited group of such entities (e.g. "halogen", "alkali metal", etc.).

Alternatively, there may be further structure diagrams, perhaps introducing new symbols for structural variables, or further citing structural variables that have already been introduced. Such structure diagrams will normally have an indication of which atom or atoms is/are attached back to the constant part.

There may be linear formulae, which can represent single entities (e.g. "OH", "COOH", "COOCH<sub>3</sub>" etc.), or include symbols for structural or multiplicative variables, or represent classes of structural entities (e.g. "C<sub>6</sub>H<sub>13</sub>").

There may be nomenclatural terms or expressions describing classes of structural entities, such as homologous series (e.g. "alkyl", "alkylcycloalkyl", "alkenyl" etc.). Frequently these are qualified by indications of the number of atoms, the degree of branching, or other factors (e.g. "straight-chain 1-6C alkyl"). Alternatively, the class described may be less well-defined (e.g. "heterocyclic ring system", "aryl" etc.).

Finally, there may be expressions describing groups in terms of their properties (e.g. "electron-withdrawing group", "photographically-useful group", "easily-hydrolysed group", "group known in the art" etc.).

In addition, all these types of expression may be further qualified by indications of position or multiplicity, or the statement that they are "substituted by" or "optionally substituted by" a further list of values. Occasionally the epithet "substituted" or "optionally substituted" may occur without any indication of the nature of the further substitution. Furthermore, certain of the alternatives listed may be indicated as preferred, possibly ranging over a hierarchy of preferability; expressions involving "preferably ... more preferably ... even more preferably ... most preferably" are not uncommon.

In some examples, two structural variables may be combined to form a ring which can be described by any of the methods given above; such combination may be a value for the two variables alternative to those given for each individually, if any. The structural entity specified as a value for the combination of the variables may consist only of the atoms added to those present in the constant part, or (more commonly) may also include those atoms of the constant part which are part of the ring formed. Occasionally two structural variables are combined to form an extra bond between the (adjacent) atoms to which they are each attached.

### 1.2.3. Conditional Expressions

Frequently, certain of the alternative values for structural and multiplicative variables are dependent upon the values of others, and this is indicated in patent specifications and abstracts in a variety of ways.

If there are several occurrences of a structural variable in the constant part, then the definition of it may specify that all its occurrences should have the same value, or different values etc. Alternatively, it may be specified that certain values for a variable are only possible when another variable has a particular value or values, or the possible values may be limited to a subset of the alternatives given originally when another variable has a particular value. There may be stipulations that a certain proportion of the occurrences of a variable should have a particular value etc. Sometimes these conditions and restrictions can become very complicated.

### 1.3. THE "MARKUSH PROBLEM"

In recent years chemical information scientists have tended to talk about the "Markush problem", and the possibilities for its solution. By this they refer to the problem of developing a computer system capable of storing and searching files of generic structures, especially those found in patents.

During the past two decades a great deal of work has been done on the development of storage and retrieval systems for specific structures, and Warr<sup>18</sup> has recently reviewed the available software. Many excellent systems have appeared, for use both with a company's files of internally-developed compounds, and with "public" databases such as the Chemical Abstracts Registry file.

Amongst the former group are the CROSSBOW system (Computerised Retrieval of Structures Based On Wiswesser)<sup>19</sup> in which structures are encoded in the Wiswesser Line Notation,<sup>20</sup> and more recently MACCS (Molecular ACCESS System)<sup>21</sup> which has sophisticated facilities for graphical input of structure diagrams. The two main systems supporting the Chemical Abstracts Registry file are CAS ONLINE<sup>22</sup> which was developed by the Chemical Abstracts Service itself, and the French Systeme DARC (Description, Acquisition, Retrieval, Correlation)<sup>23-25</sup> which is also now available for in-house use. Although these systems support limited facilities for generic structure queries, none of them, as yet, has any facilities for generic file structures.

Jackson<sup>26</sup> has outlined the essential features of an "ideal" system for generic structures in patents, and achievement of these objectives could be regarded as a solution of the "Markush problem":

1. Total recall with minimum noise.
2. Include both generic structure and specific compounds.
3. Easy to use for encoding and retrieval.

4. Automatic input with error checks.
5. Available online.
6. Abstract and structure as output.

In his paper Jackson also surveys the existing systems available, and discusses the ways in which they fall short of the ideal. Existing chemical patent documentation systems have also been reviewed by a Japanese Study Team<sup>27</sup> and in a number of other publications.<sup>28-30</sup> The storage and retrieval of Markush structures was identified as a priority area for research by the British Library's Chemical Information Review Panel, which reported in 1978.<sup>31, 32</sup>

#### 1.4. GENERIC STRUCTURE REPRESENTATIONS

An essential prerequisite for a satisfactory storage and retrieval system for generic structures is a satisfactory means of representing them for computer manipulation. A number of different forms of representation are used in existing systems and have been proposed for new systems, and these are discussed in this Section with some comments on the efficacy of the systems which use them.

Like those for specific structures, the forms of representation may broadly be divided into **ambiguous** and **unambiguous**; the former allow the same representation to stand for different structures, whereas in the latter each representation stands for only a

single structure. All operational computer storage and retrieval systems for generic structures are based on ambiguous representations of the structure, and this is one reason for the unsatisfactory performance of existing systems.

#### 1.4.1. Derwent Publications Ltd.

Derwent Publications Ltd. is a British company, owned by the Thompson Organisation, and it produces a variety of current awareness and retrospective search services, both for patents and in other areas, though patent documentation represents the major part of its business. The chemical area is well covered, and Derwent's services have been discussed recently by Kaback.<sup>33-34,</sup>  
6

In general, non-chemical patents are included in the World Patent Index (WPI), and chemically-related ones in the Central Patent Index (CPI), of which three sections (Section B on pharmaceuticals ("FARMDOC"), Section C on agrochemicals ("AGDOC") and Section E on general chemistry ("CHEMDOC")) use a complex fragmentation code, the CPI code, to represent the chemical structures, generic and specific, shown in the patent in question.

Both WPI and CPI are available for searching online via the SDC Search Service, using the ORBIT software.

The CPI code has undergone a large number of revisions during its history, which goes back to 1963 when the FARMDOC service began. It is a manually-assigned fragment code, and was originally based on the 960 punch positions available on an 80-column punched card, the cards being sorted mechanically.<sup>35</sup> The code has been substantially revised over the years, and the database made available on magnetic tape as well as punched cards, and the revisions introduced in 1982 removed the restriction to punched-card format.

Each punch position, or fragment number, represents a functional group, ring system, or other feature of chemical significance, and coding is carried out manually by highly-trained and experienced encoders; there is no automatic error checking of input. The generic structure as a whole is encoded, but this involves assigning fragment numbers for all chemical features present in the generic structure, irrespective of the logical relationships between them. Thus, in effect, all the alternative specific structures covered by a generic structure are over-coded on the same representation.

Searching is carried out by combining fragment numbers with Boolean operators, and the results are characterised by high recall and low relevance, a figure of 5% for the latter being not uncommon.<sup>26</sup>

Whilst the improvements in the code over the past two decades have been substantial, it remains less than fully satisfactory.

Up to 1977 the Pharma Dokumentation Ring (PDR), an association of European pharmaceutical companies, found it necessary to recode the generic structures from patents in Ringcode, another fragmentation code also used for Derwent's RINGDOC and Chemical Reactions Documentation Service (CRDS) services. A semi-automatic coding system, CORA,<sup>36</sup> was developed for this purpose, but in 1977 the recoding was discontinued as improvements in the CPI code had meant that Ringcode no longer gave a better retrieval performance.<sup>26</sup>

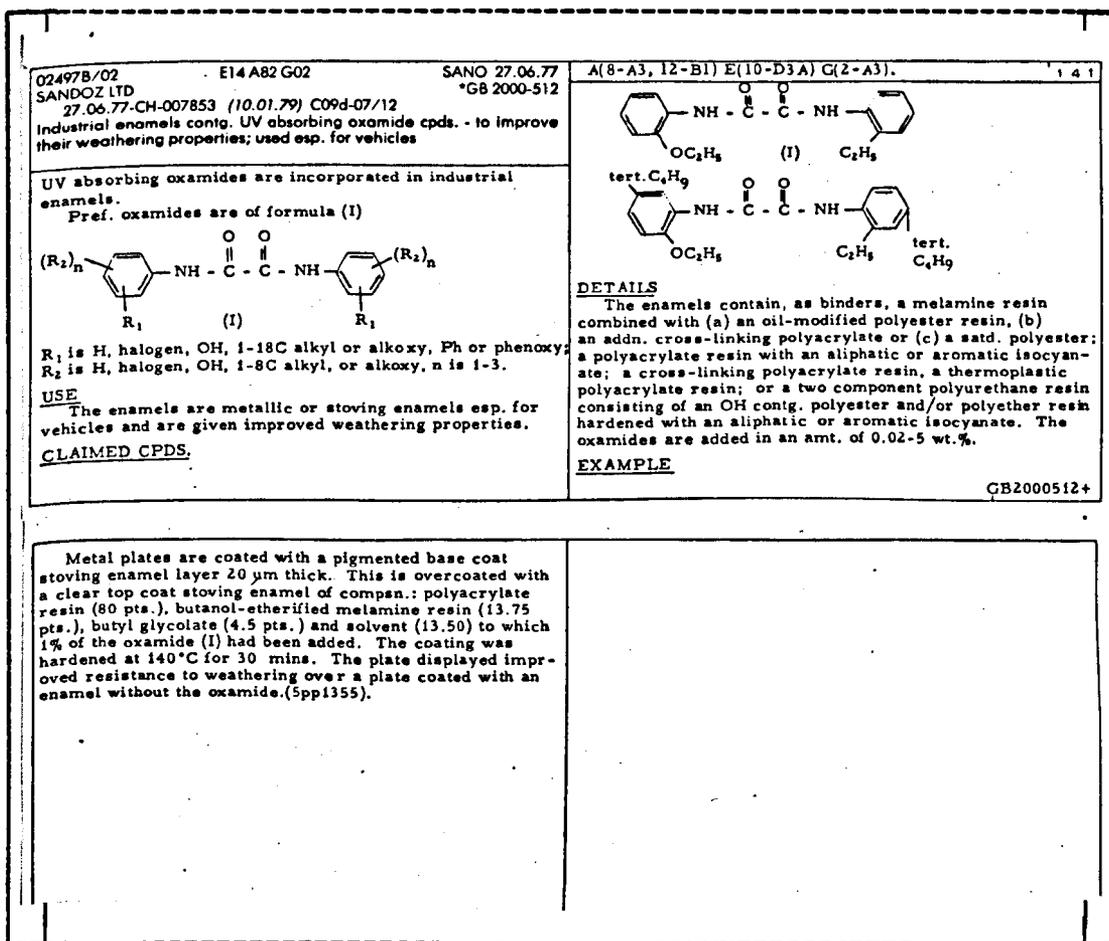


Figure 1.3: Derwent Basic Abstract for a patent.

In addition to their fragment-code indexing of chemical structures in patents, Derwent produce a compact and highly informative abstract of the patent, which was originally designed to appear on the back of the punched card used for coding. Where a generic structure appears in the patent, this is reproduced in the abstract, in which it is slightly reformatted to conform to Derwent's house style. Figure 1.3 shows the Derwent Basic Abstract for the British patent part of which was illustrated in Figure 1.2; other examples of Derwent abstracts appear in Figures 3.2 to 3.11 in Chapter 3.

#### 1.4.2. IFI/Plenum Data Co.

The patent documentation services provided by this American company have their origins in systems developed by a number of different organisations. The chemical coding system was developed by E.I. Du Pont de Nemours & Co.<sup>37-39</sup> and like Derwent's CPI code it is a manually-assigned fragment code.

Its unique aspect is that a distinction is made between fragments derived from the constant and variable parts of the structure. Figure 1.4 illustrates the assignment of such fragments for a simple generic structure, and it can be seen that those fragments deriving from either the constant or variable parts are designated possible, but only those deriving from the constant part are designated must.

In searching, the **possible** fragments are searched using positive logic, and the **must** fragments using negative logic, the latter excluding particular fragment combinations not wanted, thus improving precision.

Whilst this approach is likely to improve retrieval performance over systems such as Derwent's, which effectively use only the **possible** fragments, it does not solve the problem of indicating **possible** fragments that are mutually exclusive (e.g. halo and nitro in Figure 1.4).

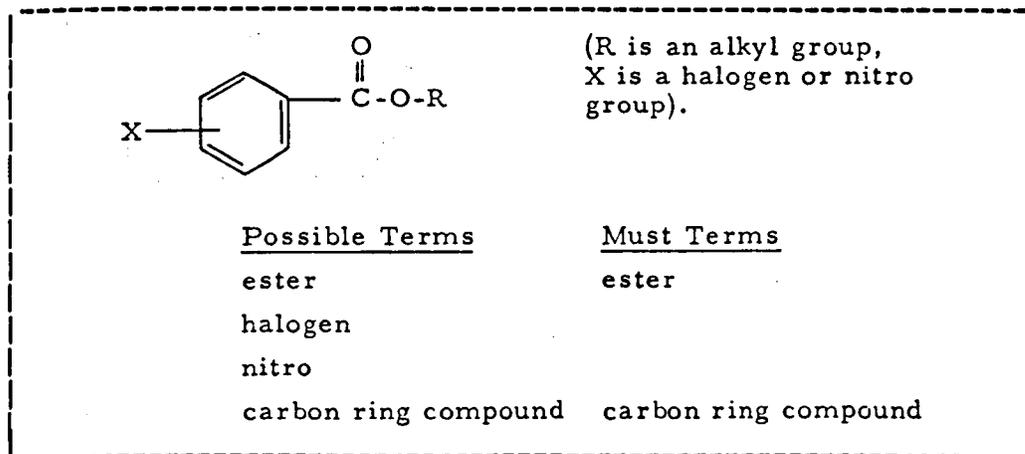


Figure 1.4: Fragments in the IFI/Plenum System

In addition to the fragment descriptors, "link" and "role" indicators are used, the former linking fragments from the same structure (where there is more than one in a patent) and the latter designating the structure as reactant, product etc. Searching can be carried out using a "weighted-term" query<sup>40</sup> in which each query term is given a "weight", retrieved documents

being those whose total score of weights exceeds a specified value.

The IFI/Plenum system, which is available online as the CLAIMS database on the Lockheed system, is restricted to United States patents, which severely limits its usefulness to patent searchers in other countries.

### 1.4.3. International Documentation in Chemistry

Internationale Dokumentationsgesellschaft für Chemie mbH (IDC) is a German company set up by a consortium of mainly German pharmaceutical companies, the principal members being BASF, Bayer and Hoechst. It is now part of the German National Information Centre for Chemistry.

So far as chemical structures are concerned, the core of the IDC system is the GREMAS (Genealogical REtrieval by MAGnetic tape Storage) code, originally developed at Hoechst.<sup>41-43</sup> This is an open-ended fragment code containing two different types of fragment, respectively called **semantic** and **syntactic** terms, and certain aspects of its design make it especially well-suited to the encoding of generic structures.<sup>44</sup> In fact Mullen<sup>45</sup> has gone as far as to claim that "the problem with Markush formulae ... [has been] solved by the GREMAS system developed by Hoechst".

The semantic terms describe the functional groups present by

means of three-letter codes, in which each successive letter indicates more precisely the nature of the group. Figure 1.5 shows some examples of the letters used to represent some common functional groups.

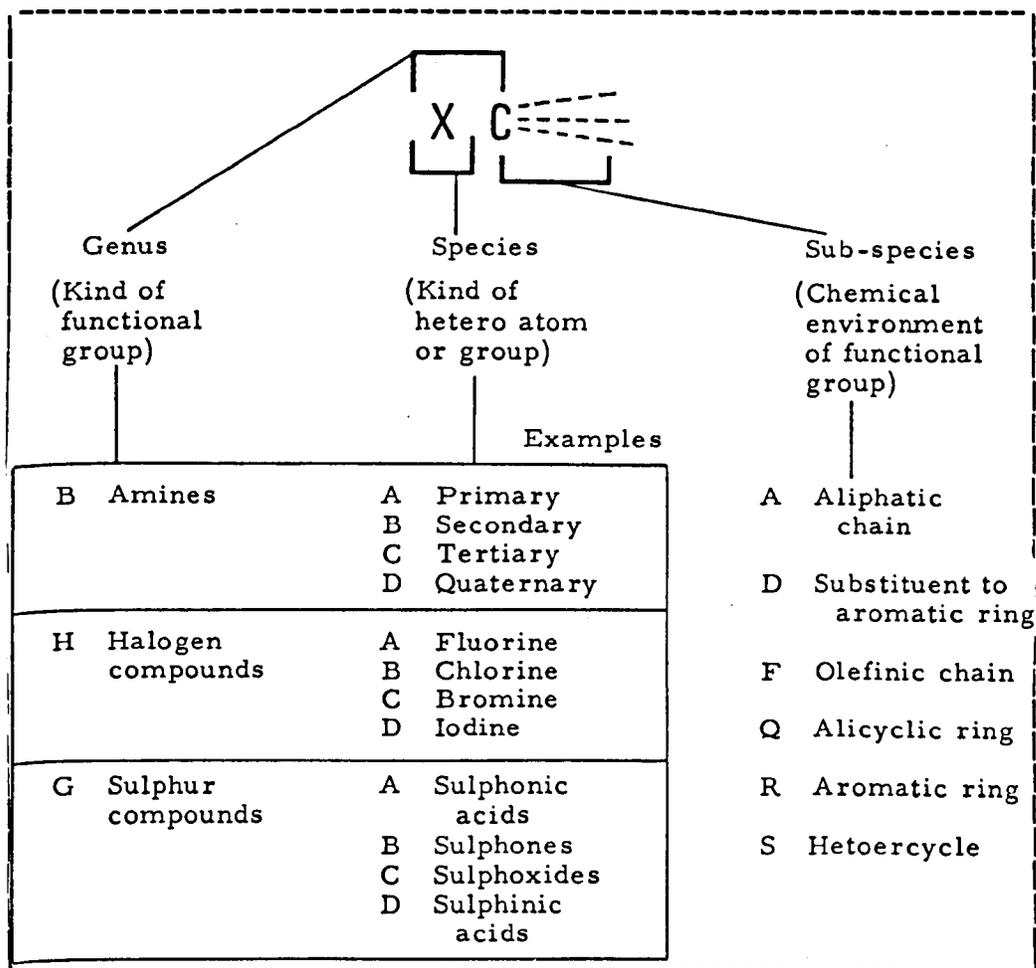


Figure 1.5: Some GREMAS semantic terms (from <sup>27</sup>)

For generic structures, the numeral 0 can be used to give a generalised fragment, e.g. HOR represents a halogen substituent on an aromatic ring, but does not specify the particular halogen.

The **syntactic** terms in the code indicate the relationships between the semantic terms, and they normally begin with a Y. Each one represents one of four **regions** of the structure, and this is indicated by the second letter: YR... for carbon chains, YS... for alicycles, YT... for aromatic rings and YU... for heterocycles. These two letters are followed by the initial letters from the semantic terms represented in the **region** in question, with the result that these terms can be of any length. Numeric locants can also be used to indicate substitution patterns on rings.

In generic structures, where there is a list of alternatives for a structural variable, the appropriate initial letters of the semantic terms are shown all together in the syntactic term, following a slash, which indicates that only one of them may be present. An example of this appears in Figure 1.6.

Specific structures can be encoded automatically in GREMAS terms, but generic structures are encoded manually.

The use of an open-ended code of syntactic descriptors in the GREMAS system, able to handle alternatives in a generic structure, makes GREMAS far more effective for storage and retrieval of generic structures than other fragment-based systems. It has, however, severe limitations in that a maximum of nine alternatives can be catered for in each of a maximum of three structural variables; its ability to handle generic expressions such as "alkyl" is also slight.

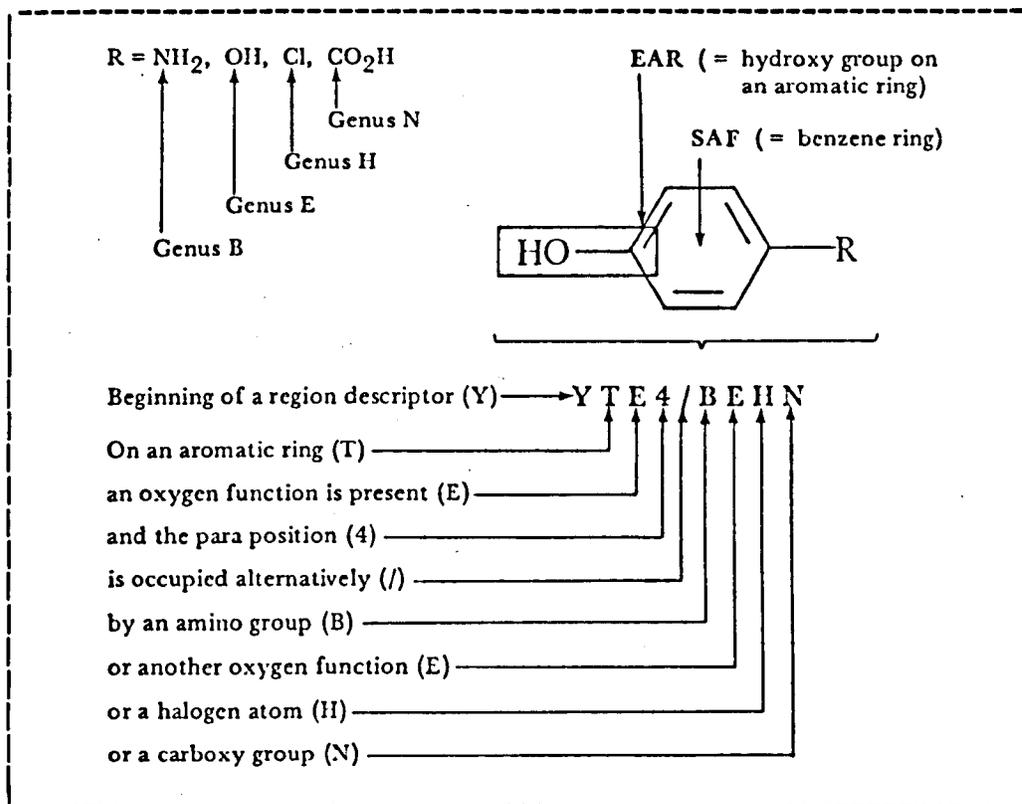


Figure 1.6: GREMAS coding for a generic structure (from <sup>27</sup>)

Silk <sup>29</sup> has pointed out that the GREMAS code possesses a far more precise search capability than other fragment codes in use, even for specific structures, since it is able to deal with specified positions of substitution on rings or chains, and to distinguish between substituents on different ring systems. The inclusion of specific structures from the Chemical Abstracts Registry file in the IDC database, along with generic structures from patents, also gives it an edge of rival systems for many types of enquiry. However, Silk also notes that the system is extremely expensive, and suggests that there could be many problems in mounting it online.

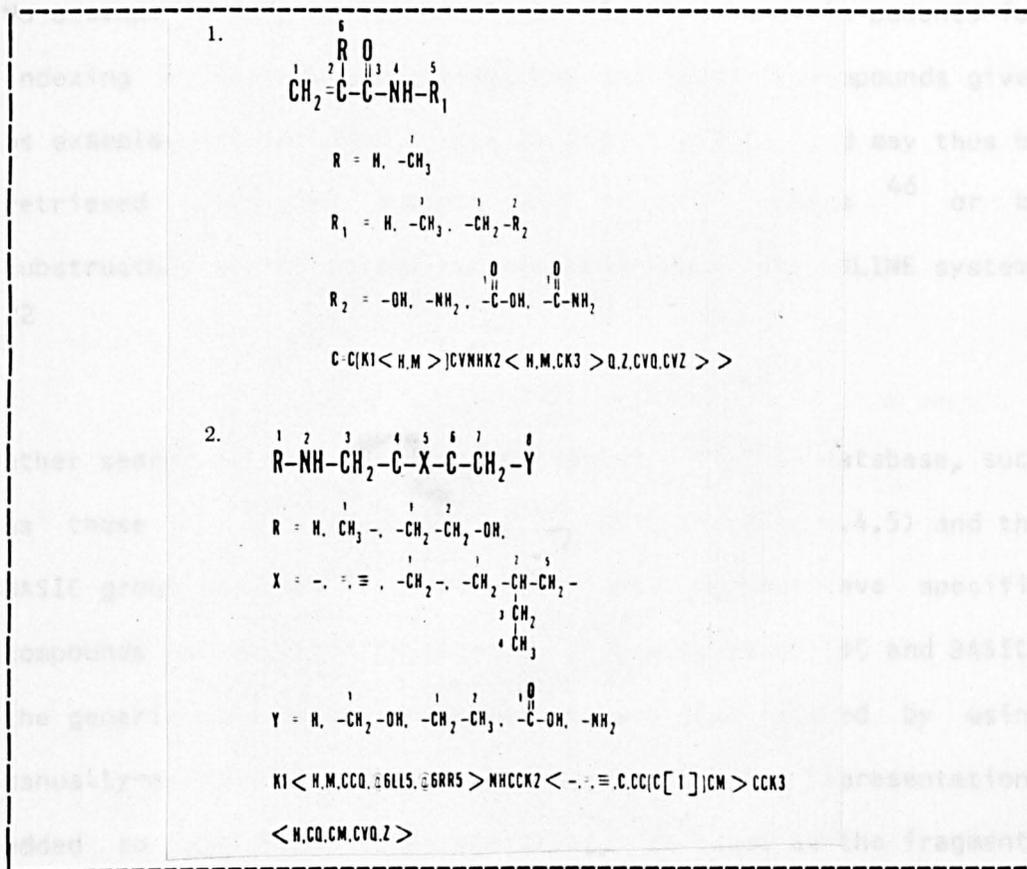


Figure 1.7: Hayward Notations for generic structures<sup>14</sup>

#### 1.4.4. Chemical Abstracts Service

Abstracts of chemical patents appear in Chemical Abstracts (CA), which also includes concordances relating basic and equivalent patents from different countries. Though in the past the coverage of patents was not as comprehensive as in other services, such as Derwent's,<sup>3</sup> it has improved recently, and CA now abstracts more than 70 000 patents annually, of which 45% are from Japan alone<sup>1</sup> resulting in serious translation difficulties.

No attempt is made to represent generic structures in patents for indexing or searching purposes, but any specific compounds given as examples are included in the CA Registry file, and may thus be retrieved using the subject and formula indexes<sup>46</sup> or by substructure search on the recently-introduced CAS ONLINE system.

22

Other search systems using the CA Registry file as database, such as those of IDC (Section 1.4.3), DARC (Section 1.4.5) and the BASIC group in Basel<sup>47</sup> are thus able to retrieve specific compounds exemplified in patents. In the cases of IDC and BASIC, the generic structures in the patents are also indexed by using manually-assigned fragment descriptors, and these representations added to the files for searching, at least at the fragment-matching level.

Whilst the effectiveness of such systems clearly depends on the relationship between the generic structure in a patent, and the specific compounds exemplified in it,<sup>48</sup> a group of searchers from ICI have suggested that even with this limitation, retrieval performance is at least comparable with that achieved by searching in a database such as Derwent's or IDC's where the generic structure is indexed by manually-assigned fragment codes.

49

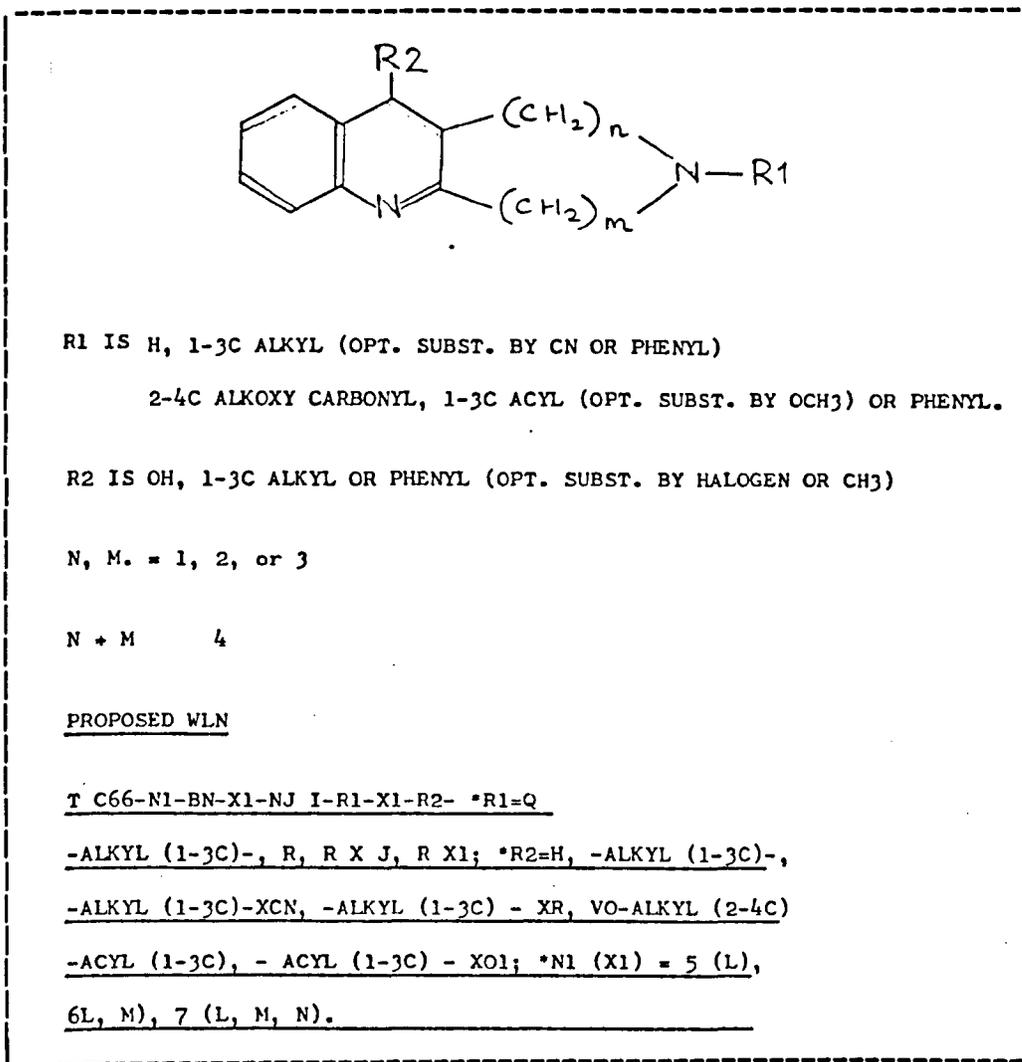


Figure 1.8: WLN representation for a generic structure<sup>26</sup>

#### 1.4.5. Système DARC

Système DARC (Description, Acquisition, Retrieval, Correlation) is a chemical substructure search system developed at the University of Paris by Dubois and others,<sup>23-25</sup> and it is available online through the French host system, Telesystèmes.

At present it permits substructure searching in the CA Registry file, and thus gives access to the specific compounds exemplified in patents. Bois and Chaumier<sup>50</sup> have compared DARC's performance to that of IDC (Section 1.4.3), and noted the latter's better coverage of patents.

Facilities for generic query structures are shortly to be implemented, though these will still only permit searching in files of specific structures. However, some comments on storage and retrieval of generic structures have appeared in publications on the DARC system<sup>51-52</sup> even if limited to the statement that "the treatment of Markush formulae has been studied in Paris by Professor Dubois", and it was recently claimed that a full generic structure search system would be ready in 1984.<sup>53</sup> No information has been forthcoming as to its capabilities or method of operation.

#### 1.4.6. Line Notations

Up until a few years ago, line notations predominated as a means of unambiguously representing specific chemical structures for machine processing, and so it was to be expected that investigations should be made into the possibility of extending such notations to handle generic structures also.

The first attempt of this sort was by the late G.M. Dyson<sup>54</sup> who showed how generic groups such as "alkyl" could be encoded in a

modification of his own IUPAC notation, along with lists of specified alternatives for structural variables. A problem area he identified was that of the definition of expressions like "cyclic carbon compound", and he suggested that a data bank might be maintained with standard notations for such expressions. This problem has also been encountered in the present work, and is discussed in Section 5.7.3.

In the mid to late 1960's work was carried out at the U.S. Patent Office on the encoding of generic structures in a form of Hayward Notation,<sup>55, 14</sup> though it was only applicable to certain types of generic structure (the so-called **determinate** structures identified by Sneed et al.<sup>14</sup> and referred to in Section 1.1.2). Figure 1.7 shows examples of the notations that resulted. Associated with this was work on search algorithms for generic structures, using a connection table representation.<sup>56-58</sup>

The dominant position of Wiswesser Line Notation (WLN)<sup>20</sup> in specific structure systems led to a number of efforts by the British software house Fraser Williams (Scientific Systems) Ltd<sup>59</sup> and others to adapt it for generic structures; an example of the rather unwieldy notations which resulted is shown in Figure 1.8.

A more promising approach was suggested by Krishnamurthy and Lynch<sup>15-16</sup> and is based on Krishnamurthy's own "ALgorithmic Wiswesser Notation" (ALWIN)<sup>60-62</sup> which is a modification of the original WLN and is designed to be amenable to automatic

generation.

A unique feature of this approach is the use of formal grammars for the representation of members of radical classes such as "alkyl". Figure 1.9 illustrates an ALWIN-based notation for a generic structure, with the associated grammar production rules.

None of these notation-based suggestions has been implemented, despite the potential advantages (discussed in Section 1.5 below) of an unambiguous representation of the generic structure. There are a number of reasons for this. In the first place the notations that result from even quite simple generic structures are generally-speaking horrendous, and a system using them could hardly be described as "user-friendly".

Secondly, many of the existing rules in line notations are designed to produce a canonical notation for a given specific structure. It is difficult to see what purpose would be served by a canonical (as opposed to merely unambiguous) notation for generic structures, whereas to ignore the canonicalisation rules altogether would result in widely-differing notations for quite similar structures.

Furthermore, the fact that many notations (WLN and ALWIN in particular) emphasise ring systems would lead to great difficulty in structures with optional rings, or with rings of variable size, on account of problems in assigning locants for substitution positions etc. Finally, the use of a line notation

for generic structure representation might severely restrict the options available for generation of fragments for a first-level screening search; it is likely that such fragments would have to be closely related to the symbols used in the notation to represent functional groups etc.

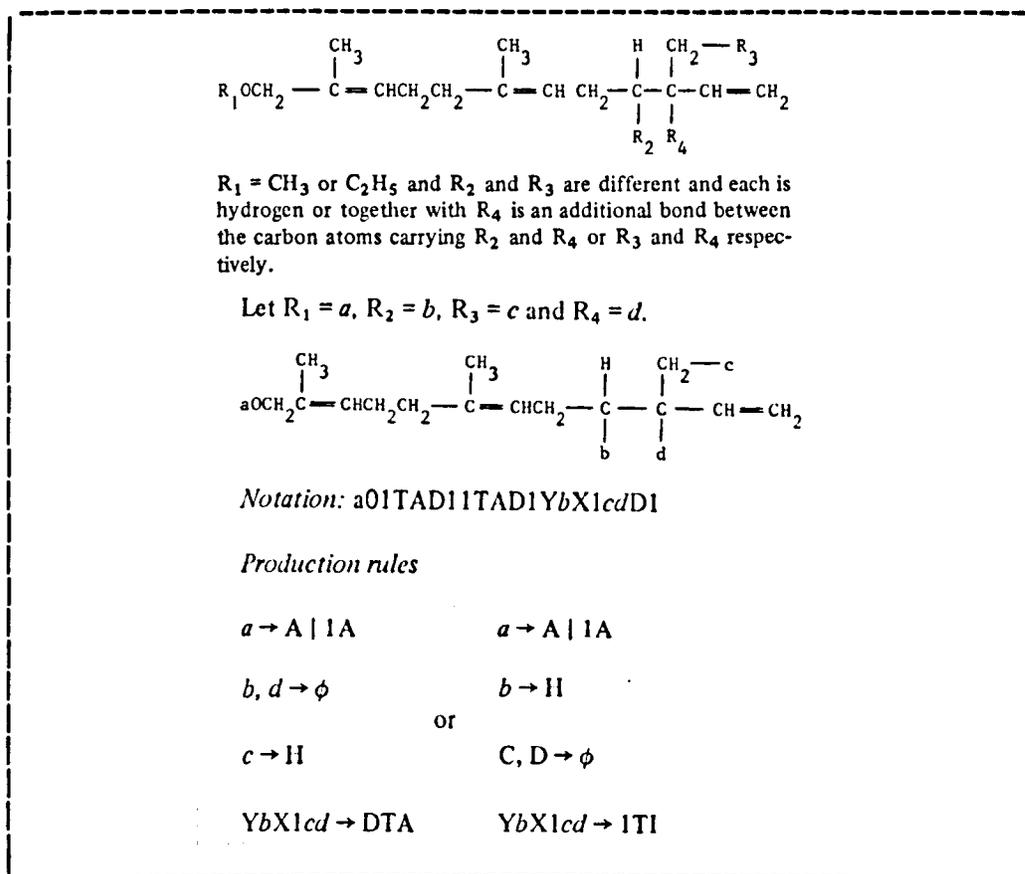


Figure 1.9: ALWIN representation for a generic structure

#### 1.4.7. The COUSIN system

Howe and Hagadone have developed an online structure storage and retrieval system at the Upjohn Company in Michigan.<sup>63-64</sup> Called COUSIN (CompOUnd Search INformation system), the system is particularly interesting on account of the extensive facilities it provides for generic query structures, though it uses a database of specific structures.

Generic queries may be input using a special notation, the  $R_k$  notation, which allows R groups to be introduced in the structure diagram for the constant part of the structure, and subsequently defined. Figure 1.10 illustrates a generic structure in  $R_k$  notation, and it can be seen that the system requires every possible attachment position for each R-group variable to be indicated in the diagram. In the definition of the R-group, each possible value is followed by the number of times it can occur in the specified positions.

The query validation program is able to check that there are no inconsistencies in the information given, and to calculate multiplicities where the user has simply specified "rest".

From the  $R_k$  notation input, the system is able to form a connection table based internal representation of the query, which is used in searching, though details of this have yet to be published.

COUSIN is not intended for use outside Upjohn, and the hardware configuration it runs on would make it extremely difficult to transport, but it is probably the most sophisticated computer representation for generic structures (albeit only query structures) currently in operation.

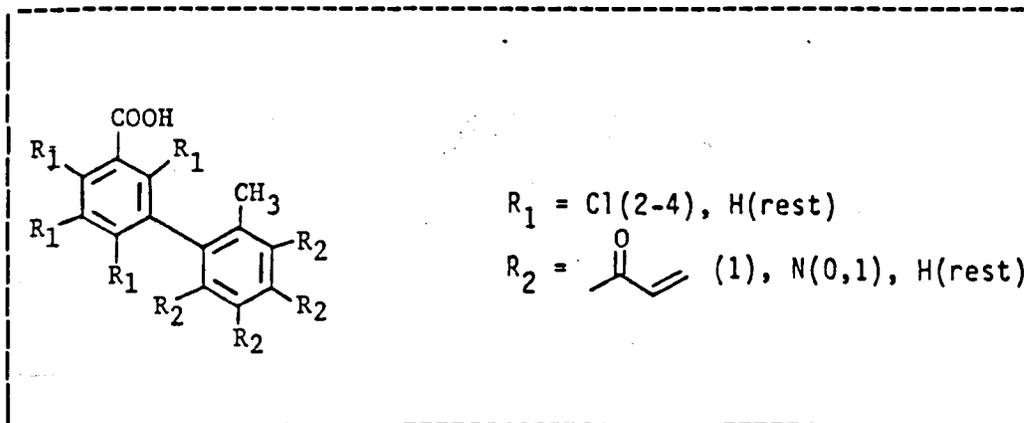


Figure 1.10:  $R_k$  notation for a generic structure<sup>64</sup>

### 1.5. REQUIREMENTS FOR A SEARCH SYSTEM REPRESENTATION

The various forms of generic structure representation described in the last section are all unsatisfactory for one reason or another, and the work described in this Thesis has had as its aim the development of a more effective representation, allowing a closer approach to Jackson's concept of an "ideal" generic structure information system (Section 1.3). This has led to the idea of a number of different representations for use at different stages of such a system.

The conventional arrangement of storage and retrieval systems for specific structures has involved a number of stages in a search. Lynch<sup>65</sup> has discussed the need for a first-level "screening" search to remove from consideration those structures in the database which, by virtue of their lack of some feature present in the query, cannot possibly satisfy the query. When the file to be searched has been thus reduced, computationally more expensive procedures can be used to search those structures which remain candidates. A variety of different "screens" have been used for first-level searching, including molecular formulae and various fragment-code representations, often implemented as bit-screens.

The present work has envisaged an analogous approach to generic structure searching, and Figure 1.11 illustrates the overall process intended. An input notation, called GENSAL (GENERIC Structure Language), has been designed for the unambiguous description of generic structures in a form which is intelligible to a chemist or patent agent, yet sufficiently well formalised to permit automatic analysis by computer. GENSAL is intended to be the representation used for input both of file structures from patents and of query structures, and it is described in Chapter 3. It is a formal language, analogous to a computer programming language, and Chapter 2 reviews briefly the theory of such languages.

The GENSAL representation input to the computer will be used to generate an internal representation of the structure, and this is described in Chapter 4. It is based on connection tables<sup>66</sup> and

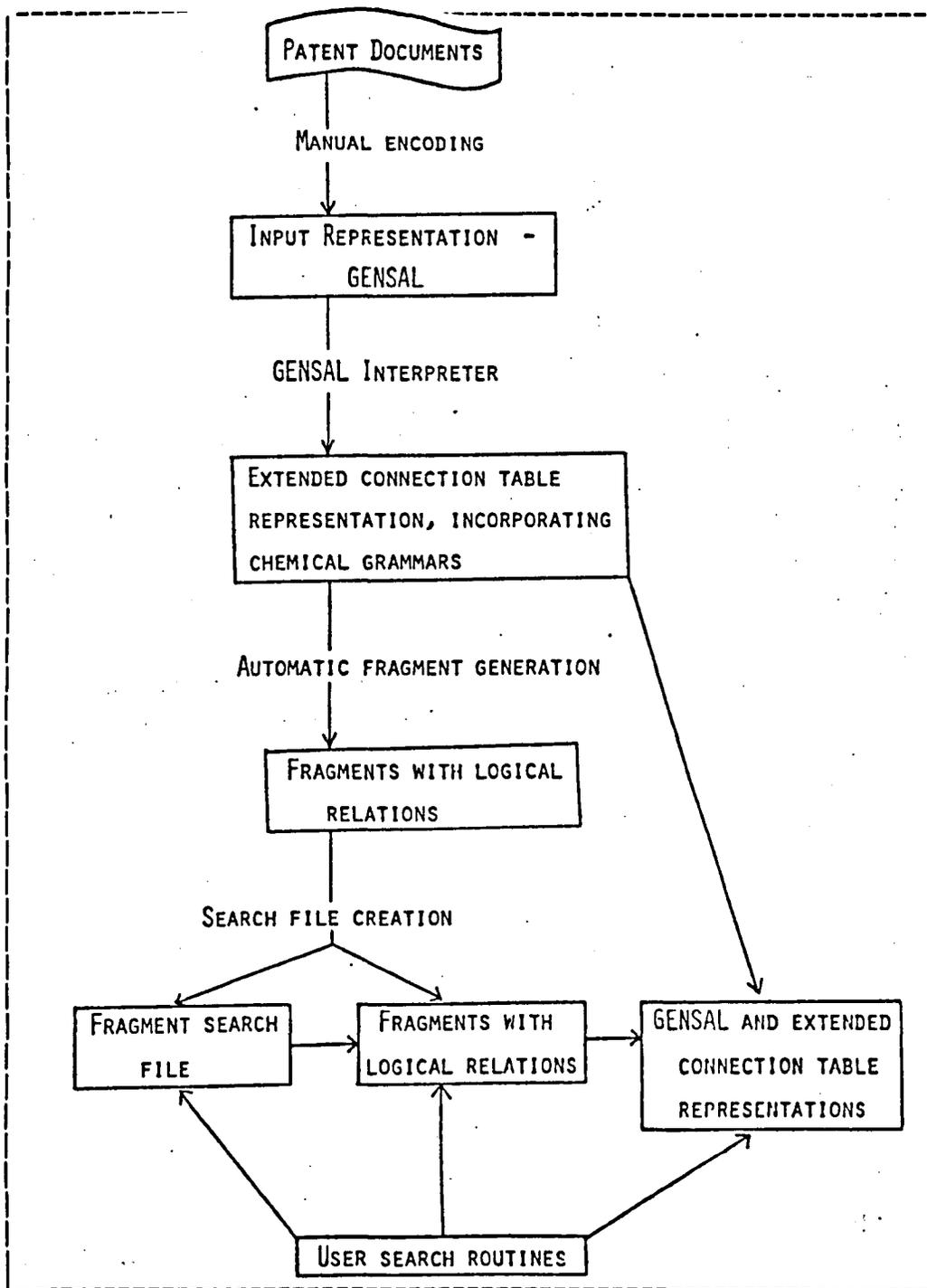


Figure 1.11: Overall process for generic structure system

Like GENSAL is an unambiguous representation of the generic structure; it is intended to be transparent to the user. The interpreter program which performs the conversion from GENSAL to

the internal representation is described in Chapter 5.

The internal representation is envisaged as the basis for searching. Ultimately, it should be possible to perform an atom-by-atom match between query and file structures in the internal representation, but such a match is likely to be extremely expensive computationally - much more so than in specific structure search systems, on account of the possibility of alternatives at various points.

Thus it is expected that there will probably be at least two fragment-based screening searches to reduce the file of candidate database structures. A number of different types of fragment can be generated from the internal representation, and algorithms for such fragment generation, and the use of fragments in different search representations are discussed by Welford.<sup>67</sup>

## CHAPTER 2

### FORMAL LANGUAGES

"I conceive you may use any language you choose  
to indulge in without impropriety."

W.S. Gilbert

The mathematical theory of languages has been extensively developed over the past quarter of a century, and has been the subject of several textbooks.<sup>68-73</sup> This Chapter gives an outline of those aspects of the theory of formal languages, and the means of parsing them, which have been built upon in the design of the GENSAL language described in Chapter 3, and in the programming of its interpreter, described in Chapter 5. It will do this with particular reference to computer programming languages, which are the most commonly encountered class of artificial formal

Languages, and in this context will discuss the choice of programming language for the software development described in Chapter 5. The Chapter also considers the use of artificial formal languages in information work, and in particular in chemical information.

No attempt will be made to give a comprehensive review of the subject of formal language theory, as many excellent such reviews exist, and will be referred to, and as far as possible the more mathematical aspects of the area will be avoided.

### 2.1. DEFINITION AND CLASSIFICATION OF FORMAL LANGUAGES

The earliest work on the mathematical theory of languages, which was done in the late 1950's, is largely due to Noam Chomsky,<sup>74-78</sup> who was attempting to find a means of modelling natural languages such as English. His aim was to understand the mechanism by which it is possible to comprehend sentences never heard before, and to produce completely novel, but grammatically correct, sentences.

For the purpose of his analysis Chomsky considered a language as being a set (finite or infinite) of sentences, each finite in length and constructed by concatenation out of a finite set of elements. These elements are termed the "terminal symbols" of the language, and might, in the case of English, be identified with the set of valid English words.

The grammar of a language he considered as a means for generating sentences in such a language. The grammar will generate all possible grammatically correct sentences in the language, but no others. It specifies the symbols of language, and includes a set of rules, sometimes called "productions", or "rewriting rules" which specify the replacement of one group of symbols by another during the generation of a sentence: a grammatically correct, or "well-formed" sentence in a language is one that can be generated by the grammar.

Whilst a given grammar is only able to generate sentences in a single language, several different grammars may all generate the same language - such grammars are said to be equivalent.

Put more mathematically, a grammar  $G$  may be represented as a "4-tuple":

$$G = ( V_N, V_T, P, S )$$

$V_N$  is the set of "non-terminal symbols" or "variables" (descriptive terms or "metasymbols" representing elements of the sentence), and  $V_T$  is the set of "terminal symbols". Both  $V_N$  and  $V_T$  are called "alphabets", and they are disjoint. Their union is symbolised  $V$ .

Strings, or "sentences", can be constructed over an alphabet, and consist of concatenated sequences of elements of the alphabet, of arbitrary length. The set of sentences over an alphabet  $V$  is

symbolised  $V^*$ , and may include the null string (which is of zero length). The set of sentences over  $V$ , excluding the null string, is symbolised  $V^+$ .

$P$  is a set of "productions" or "replacement rules", which are of the form

$$\alpha \rightarrow \beta$$

where  $\alpha$  is a string in  $V^+$  and  $\beta$  a string in  $V^*$ .

If a production in  $P$  can be used to rewrite a string  $\alpha_1$  as another string  $\alpha_2$  then it is said that  $\alpha_1$  directly derives  $\alpha_2$  in grammar  $G$ . If the application of a series of productions in  $P$  enable  $\alpha_1$  to be rewritten as  $\alpha_m$  then it is said that  $\alpha_1$  derives  $\alpha_m$  in grammar  $G$ .

The grammar  $G$  is said to generate a language  $L(G)$ , which consists of the set of sentences over  $V_T$  (i.e. elements of  $V_T^*$ ). However, only certain of the sentences in  $V_T^*$  are grammatically correct ("well-formed"), and these are those of them that can be derived from  $S$  (which is a distinguished member of  $V_N$  called the "sentence symbol" or "start symbol") in grammar  $G$ .

2.1.1. The Chomsky Hierarchy

grammars, and hence the languages which they generate, have been classified by Chomsky <sup>74</sup> by imposing successively tighter restrictions on the form of the production rules in P. The above grammar, in which no restrictions are imposed is called a Type 0, or unrestricted grammar, and the languages it generates are called the recursively enumerable languages.

A Type 1 grammar is obtained if it is required that the number of symbols on the right hand side of each production should be greater than or equal to the number of symbols on the left hand side. An alternative, and equivalent restriction is that the rules in P should be of the form

$$\alpha_1 A \alpha_2 \rightarrow \alpha_1 \beta \alpha_2$$

where  $\alpha_1$  and  $\alpha_2$  are in  $V^*$ ,  $\beta$  is in  $V^+$ , and A is in  $V_N$ . This form of the restriction leads to the name context sensitive for this type of grammar, as it allows A to be replaced by  $\beta$  when it occurs in the context of  $\alpha_1$  and  $\alpha_2$ .

In Type 2 grammars, the left hand side of the production must be a single non-terminal symbol. The productions are therefore of the form

$$A \rightarrow \beta$$

where  $A$  is an element of  $V_N$ , and  $\beta$  a string in  $V^*$ . Since  $A$  may be replaced by  $\beta$  independently of the context in which it occurs, this type of grammar is called context free.

The most restricted type of grammar, Type 3, requires that all productions are of the form

$$A \rightarrow a B$$

or

$$A \rightarrow B$$

where  $A$  and  $B$  are members of  $V_N$  and  $a$  is a member of  $V_T$ . Type 3 grammars are called regular grammars.

It is clear that these increasingly severe restrictions on the form of the productions mean that the types of grammar are arranged in a hierarchy: every Type 3 grammar is also Type 2, every Type 2 grammar is Type 1, and every Type 1 grammar is Type 0. This is sometimes referred to as the Chomsky hierarchy.

Many important and interesting properties can be shown for all these grammar types, 68-70, 72, 74, 77, 79-81 but detailed coverage of them is beyond the scope of this Thesis.

2.2. PARSING OF CONTEXT-FREE LANGUAGES

The relative simplicity of context-free grammars has allowed considerable progress to be made in the automatic syntactic analysis (parsing) of sentences in the languages generated by them, whereas such analysis has proved highly intractable for context-sensitive and unrestricted grammars.

Unfortunately, despite initial hopes, context-free grammars have not proved adequate for the description of natural languages, but they have been extremely useful in the definition of artificial languages, in particular, programming languages.

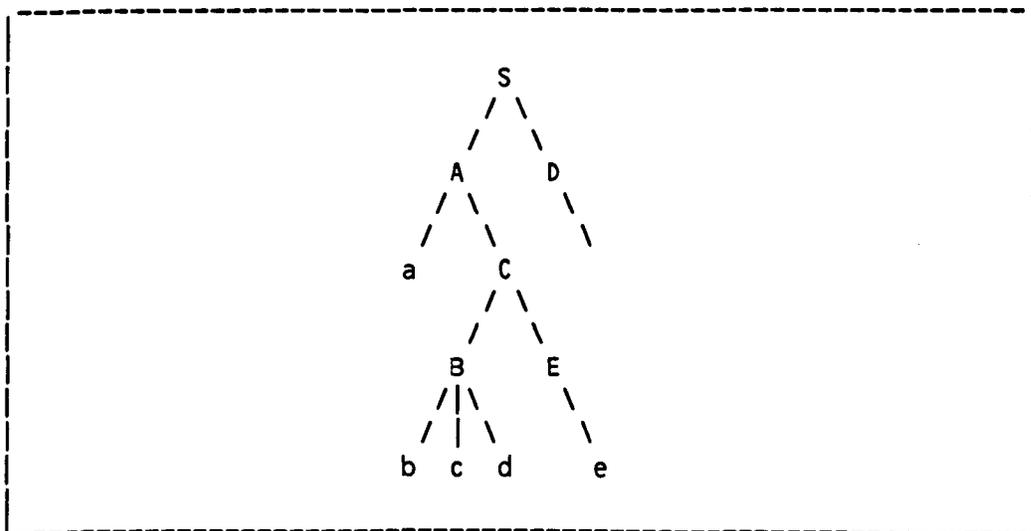


Figure 2.1: Derivation tree for the sentence  $abcde$  in  $L(G_1)$ .

A sentence in a context-free language can be analysed in terms of a grammar which generates it using a derivation tree or parse diagram in which the root vertex of the tree is  $S$ , its leaves are all elements of  $V_T$ , and its interior vertices are all elements of

$V_N$ . Consider the language  $L(G_1)$  generated by grammar  $G_1 = (V_T, V_N, P, S)$  where

$$V_T = \{ a, b, c, d, e \}$$

$$V_N = \{ A, B, C, D, E \}$$

$$P = \{ S \rightarrow AD$$

$$A \rightarrow aC$$

$$B \rightarrow bcd$$

$$C \rightarrow BE$$

$$D \rightarrow$$

$$E \rightarrow e \}$$

The derivation tree for the sentence  $abcde$  is shown in Figure 2.1.

It will be seen that the symbols of the sentence can be followed around the leaves of the tree from left to right, and that each branch of the tree corresponds to an appropriate production in  $P$ .

If a sentence has more than one derivation tree, then it is ambiguous -- grammars which generate only sentences with unique derivation trees are unambiguous grammars, whereas other grammars are ambiguous. Since it is likely to be essential for an artificial language to have an unambiguous grammar, this point is very important.

Even when there is a unique derivation tree, it is possible to obtain it by applying the productions in different sequences. The sequence may, however, be standardised by stipulating that at

each derivation the rightmost non-terminal symbol should be replaced in accordance with an appropriate production. The derivation sequence corresponding to Figure 2.1 is therefore:

$S \rightarrow AD \rightarrow A \rightarrow aC \rightarrow aBE \rightarrow aBe \rightarrow abcde$

### 2.2.1. LR Parsing

In parsing a sentence, it is necessary to start from the string of terminal symbols (i.e. the leaves of the tree), and to reconstruct the derivation tree, leading eventually back to  $S$ , its root. This is a process of successively reducing substrings of terminal and non-terminal symbols in accordance with the productions.

At each step in the parse, the derivation required is, according to the above stipulation, one in which the rightmost non-terminal symbol is replaced. This corresponds, in the parsing direction, to reducing the leftmost set of adjacent leaves of the tree (which will only all be in  $V_T$  at the start of the parse) that form a complete branch. Since the parsing thus operates from left to right along the original string of terminal symbols, it is called a left-right parse.

Knuth<sup>82</sup> has defined a subclass of context-free grammars, called LR(k) grammars, for which this parsing method will work, looking ahead a maximum of  $k$  symbols to identify with certainty each

production in the derivation. It may also be shown that all LR(k) grammars are unambiguous, and that there is an algorithm to determine whether or not a context-free grammar is LR(k) for a given k.

The derivation tree is reconstructed from the bottom upwards, and LR(k) grammars are therefore sometimes called bottom-up grammars, and the parsing method bottom-up parsing.

### 2.2.2. LL Parsing

Lewis and Stearns<sup>83</sup> have defined another subclass of context-free grammars, the LL(k) grammars, which allow an even simpler approach to the parsing of sentences. Each production in the derivation can be identified with certainty by inspecting the sentence from its beginning (left) end to the k-th symbol beyond the beginning of the production.

In this type of parse, the derivation tree is being reconstructed from the top downwards, and hence LL(k) grammars are called top-down grammars, and the parsing method top-down parsing, or parsing by recursive descent.

At the start of a parse based on an LL(k) grammar, it is assumed that a production having S as its left-hand side is required. Which of the various productions in P having S as left-hand side is appropriate can, for an LL(k) grammar, be determined by

Looking at a maximum of  $k$  symbols.

Rosenkrantz and Stearns<sup>84</sup> have shown that it is possible to determine if a grammar is  $LL(k)$  for a given  $k$ , that all  $LL(k)$  grammars are unambiguous, and that the  $LL(k)$  grammars are a subset of the  $LR(k)$  grammars. In addition they have shown that provided there are no productions with an empty string as right-hand side, it is possible to construct for a language generated by an  $LL(k)$  grammar an equivalent  $LL(k)$  grammar in Greibach Normal Form (i.e. where the right-hand side of each production starts with a terminal symbol).<sup>85</sup> Furthermore, if every production with a given non-terminal as its left-hand side has a different terminal as the first symbol on its right-hand side, then the grammar is  $LL(1)$ , and is a member of the class of simple deterministic grammars described by Koranjak and Hopcroft.<sup>86</sup> No look-ahead is required in the parsing of sentences in languages generated by such grammars, and the production involved can be determined at each step simply by examining the next symbol in the sentence. Some properties of deterministic context-free languages have been discussed by Ginsburg and Greibach.<sup>79</sup>

### 2.2.3. Top-Down vs. Bottom-Up Parsing

The properties of  $LL(1)$  and  $LR(1)$  grammars have recently been compared by Beatty<sup>87</sup> and Knuth<sup>88</sup> has compared the particular advantages of top-down and bottom up parsers. Because all  $LL(k)$  grammars are also  $LR(k)$ , any language that can be parsed top-down

can also be parsed bottom-up, but the reverse is not true, making bottom-up parsing more generally applicable. The chief advantage of top-down analysis is that it is known which production is being used after examining only  $k$  terminal symbols, and this enables some degree of prediction on the part of the parser as to which symbols will be encountered next. This is especially helpful for the semantic analysis of the sentence, and the design of modern programming languages has taken particular note of the advantages of LL( $k$ ) grammars with a low value for  $k$ .

More detailed discussion of parsing methods for context-free languages may be found in a number of textbooks and reviews.<sup>71, 80, 81, 89</sup>

### 2.3. PROGRAMMING LANGUAGES

The earliest developments of high-level computer programming languages took place in isolation from the work of Chomsky and others on formal languages, and a comprehensive survey of their history has been given by Sammet.<sup>90</sup> Fortran was the first high-level language to gain wide acceptance, and it is still the most commonly-used language for scientific applications.

Attempts have been made to formalise the grammar of Fortran<sup>91</sup> but on account of the rigid field format for its statements, and the numerous minor restrictions on various constructs, these have been of limited success. The language was designed for speed of

execution, rather than simplicity of syntax analysis.

The ALGO<sup>r</sup>ithmic Language Algol 60 was designed by an international committee in the late 1950's and early 1960's<sup>92</sup>,<sup>93</sup> and marked a turning point in programming language design. Its importance lies more in the manner of its definition, which has had a major influence on the design and definition of more recent languages including Algol 68,<sup>94</sup> Pascal<sup>95</sup> and Ada,<sup>96</sup> than in its actual use, which has been comparatively limited, at least so far as computer implementations are concerned, though it is the standard publication language for algorithms.

### 2.3.1. Syntax Specification

The original definition of Algol 60<sup>92</sup> first introduced the so-called Backus-Naur metalanguage for the formal specification of its syntax. An example of a grammatical rule of Algol 60 expression in Backus-Naur Form (BNF)<sup>97</sup> is

```
<conditional statement> ::= if <boolean expression> then
<statement> else <statement> | if <boolean expression> then
<statement>
```

This defines the syntactic category "conditional statement" as being one of two alternatives: either the word "if" followed by a

"boolean expression" followed by the word "then" followed by a "statement" followed by the word "else" followed by another "statement", or alternatively the word "if" followed by a "boolean expression", followed by the word "then" followed by a "statement". The syntactic categories "boolean expression" and "statement" are defined by other rules in the grammar.

The grammatical rules in a BNF grammar are expressed in a "metalanguage", which uses certain symbols that do not occur in the language being defined. The symbol ::= means "is defined to be", and | means "or". Angle brackets are used to enclose the names of syntactic categories, which thus themselves form symbols of the metalanguage. The words not so enclosed (if, then, etc.) are of course, the actual symbols of Algol 60.

In 1962 Ginsburg and Rice<sup>98</sup> proved that Algol-like languages defined using a BNF metalanguage are equivalent to the context-free languages (Type 2) defined by Chomsky which allowed the rigorous mathematical properties of context-free languages to be applied to Algol and other programming languages. In the BNF metalanguage, the syntactic categories can be identified with the non-terminal symbols of Chomsky Type 2 grammars and the actual symbols of the language being defined with the terminal symbols. For a programming language such as Algol 60, the start symbol is identified with the syntactic category "program". The various alternatives separated by the | symbol, correspond to the different productions having the same left-hand side.

The only difference between grammar specification using BNF or a 4-tuple is in the form of its representation. Other methods for syntax specification have also been suggested, including a tabular format<sup>99</sup> and the use of "syntax diagrams".<sup>100</sup> These latter represent the BNF rules in diagrammatic form, with separate branches for each alternative, and it is normally possible to combine several BNF rules into a single diagram. Wirth<sup>101</sup> has also proposed an extended BNF formalism.

### 2.3.2. Syntactic Analysis

The purpose of a high-level programming language is to allow a programmer to give instructions to a computer in a form which remains reasonably intelligible to himself, or to another programmer. Before the computer can actually execute the instructions, however, it must convert them into a form more closely related to its own internal architecture, and this process of conversion is called compilation.<sup>80, 81</sup> Three principal operations are involved in compilation: lexical analysis (in which the string of characters forming the program in the high level source language is split up to identify the separate tokens or terminal symbols of the language), syntax analysis (in which the grammatical relationships between the tokens are identified, in accordance with the rules of the grammar) and code generation (in which the machine level object language is generated).

The process of syntax analysis is based on the same principles of parsing as are described in general terms in Section 2.2 above, and is obviously much simpler for a programming language with an appropriately simple grammar. In the design of Algol 60 and the languages based on it particular attention has been paid to the need for simplicity of syntax analysis. Not only does this simplify the complexity of the program required to perform the syntax analysis, but a simple grammar also makes it much easier for the programmer to write elegant and error-free programs.

Irons<sup>102</sup> described a bottom-up syntax analyser for Algol 60 in 1961, but the first top-down syntax analyzer for a programming language was written for Cobol, and described by Conway<sup>103</sup> in 1963. Numerous textbooks and reviews consider the problems of compiler writing and syntax analysis for a variety of programming languages. 71, 80, 81, 104-110

### 2.3.3. The Pascal Language

Pascal is a high level language based on Algol 60, and was designed by Nicklaus Wirth, who published the first description of it in 1970,<sup>111</sup> with a revised version in 1975.<sup>95</sup> A committee of the International Standards Organisation convened by A.M. Addyman has drafted a Standard definition for the language, which has been published for comments.<sup>112, 113</sup> The language has become extremely popular, particularly in academic circles, and has been the subject of many textbooks.<sup>114-117</sup>

Pascal was designed especially for compilation using top-down (or "recursive descent") syntax analysis, and Wirth described the first compiler in 1971.<sup>118</sup> The first version of this was written in Pascal itself, and manually translated into a lower level language. Each subsequent version of the compiler could then also be written in Pascal, and compiled by the previous version, a procedure known as "boot-strapping".

Pascal has been enthusiastically promoted by many authors<sup>119-122</sup> and possibly partly as a result of this has also attracted considerable criticism,<sup>123-128</sup> some of it quite vitriolic.<sup>123-126</sup> Other authors, whilst generally welcoming the language, have made suggestions for its enhancement,<sup>129, 130</sup> and Wirth himself has published his own retrospective assessment.<sup>131</sup>

In his paper, Wirth discusses the advantages Pascal has for the writing of reliable software, as the design of Pascal permits a great deal of checking on the self-consistency of the program to be done by the compiler, and thus a high proportion of program errors can be detected before execution begins. Its highly structured design also makes it suitable as a teaching language, and this has been its principal area of application to date. Conradi<sup>124</sup> has however pointed out some of its disadvantages as a systems programming language, particularly with its lack of flexibility in matters such as the absence of dynamic arrays (Pascal, unlike for example Algol 68, requires that array bounds be specified at compile time), though Wirth's paper points out that it is precisely these limitations on flexibility that give

Pascal its enhanced security.

Despite its acknowledged weaknesses, Welsh, Sneeringer and Hoare have expressed

"the belief that Pascal is at the present time the best language in the public domain for the purposes of systems programming and software implementation".<sup>126</sup>

#### 2.3.4. The Ada Language

Pascal was used as the basis for all the tenders to the U.S. Department of Defence for the design of a new programming language to be used for all their software development.<sup>132</sup> However, the selected language, Ada,<sup>96</sup> has led to even fiercer controversy than Pascal.<sup>133-136</sup> Much of the criticism has attacked the increased flexibility of Ada over Pascal, with many additional features not present in the older language, which, it critics claim, make it less secure, and programs written in it unreliable. In view of the likely military applications of Ada, echoes of this discussion have reached a public forum.<sup>137</sup>

### 2.3.5. Choice of Language for Software Development

In choosing a programming language for the practical work described in Chapter 5 of this Thesis a number of factors were considered. A modern, structured language was required, with good program readability and portability, since the work is of substantial interest to the chemical and patent documentation industries. In addition, the programs required would operate interactively, and would therefore need to be developed on the Sheffield University Prime computer system, which restricted the choice of language to those for which Prime compilers were available.

Whilst Fortran would have provided the greatest portability, it was felt that it was insufficiently well-structured, the same reservation applying to Basic. Implementations of both Algol 68 and Pascal were available, but only the latter was actively supported by Computing Services staff, and was therefore the language chosen.

Initially a compiler developed at the University of Hull was used, but it was later replaced by a much more powerful one written by staff of Sheffield University Computing Services,<sup>138</sup> which generates segmented object code, and allows much bigger programs and easier interface with routines in other languages, and contains facilities for separate compilation of Procedures and Functions. It was also found that the easy availability of the compiler's writers was extremely useful on encountering

problems in software development; none of these advantages would have been available with Algol 68.

Nevertheless, a number of disadvantages were encountered with Pascal, of which the most serious was in the use of external files, particularly as Pascal does not permit programs to append data to files that already exist, and neither does it implement direct-access files.

#### 2.4. FORMAL LANGUAGE SEMANTICS

Chomsky<sup>78</sup> has pointed out that there may be sentences in a language which, whilst being grammatically correct, make no sense. An English example he gives is the sentence

"Colourless green ideas sleep furiously."

Similar problems may be encountered in programming and other formal languages, and though methods for specifying the syntax of a language (at least for certain classes of language) are now well-established, comparatively little success has so far been achieved in formally specifying the semantics of languages.

Several approaches have been used, and have been reviewed by a number of authors.<sup>139-141</sup> Hoare and Wirth have attempted<sup>142</sup> to define the semantics of Pascal rigidly, using an axiom-based method developed by Hoare.<sup>143</sup>

The division between the syntax and semantics of a programming language is not a sharp one, and not all authors agree on where it lies. Essentially, syntax is concerned only with those matters that can generally be defined with reference to the sequence of symbols in sentences of the language; semantics is concerned with everything else.

Wirth <sup>109</sup> has pointed out that even where the syntax of formal languages is context-free, its semantics may be context-dependent. In programming languages, semantics is concerned with such matters as type compatibility in expressions and assignments. For example the Pascal expression

5 + 'B'

is valid syntactically but not semantically as the **integer** constant 5 is not of the same type as the **char** constant 'B'. had this expression been the controlling expression in a **while** loop, then additionally its resultant type would have had to be **boolean**: this exemplifies the context-dependency of formal language semantics - even when the syntax is context-free - which is one of the difficulties in the way of the achievement of formal semantics.

Ultimately, it is the implementation of a programming language in a compiler that defines its semantics; in written descriptions of the language the semantics is normally defined informally. In any case, it is often useful to leave certain aspects of the

semantics implementation-dependent as the most appropriate way of implementing them may depend on the machine architecture in question. The type `char` and the value of `maxint` are two aspects of Pascal deliberately left undefined for this reason.

### 2.5. INTERACTIVE LANGUAGES

For most programming languages, the operation of compilation requires no interaction with the programmer, and is often carried out in batch mode. Once compiled, the object program produced by the compiler can be executed repeatedly on different data, without recompilation.

In such a system the compiler reports any error (syntactic or semantic) that it encounters, and then attempts to continue to process the source program and to report any further errors. Obviously it is no longer practicable to continue to generate object code. This has the advantage that the programmer has all the errors in his program reported together, and can correct them all before attempting to recompile it, but has the disadvantage that the compiler may not successfully recover from an error it encounters, and may then report large numbers of spurious errors. During the present work, the author had over two hundred errors reported after a compilation, all of which were corrected by the addition of a single semicolon near the top of the program.

For certain applications, compilation may take place

interactively. The programmer types his program into the computer line by line, with the compiler reporting each error as soon as it is encountered, and the programmer correcting it immediately. Languages compiled in this way are usually specially designed for the purpose, and have been discussed by Kupka and Wilsing.<sup>144</sup> The interactive compilers used for such languages are normally called "interpreters" to distinguish them from batch-mode compilers, and the special problems of writing them have been discussed by Brown.<sup>105</sup>

These authors point out that systems based on interactive compilation actually require three different languages - the programming language itself, a Command language which controls such matters as the saving of completed programs, execution etc., and an Edit language which allows interactive editing of the program. This latter is especially useful for correcting errors which are only detected by the compiler some time after they have occurred.

Both the Edit and Command languages are normally very simple, each "sentence" consisting only of a single terminal symbol (e.g. a Command) followed by one or two arguments such as a filename or a line number. Their syntactic analysis is trivial.

The language most commonly implemented in this fashion is Basic<sup>145, 146</sup> though the approach has also been applied, at least for teaching purposes, to Fortran,<sup>147</sup> Algol<sup>148</sup> and Pascal,<sup>149, 150</sup> in the latter cases only a subset of the language being

implemented. In the case of higher-level programming languages, certain problems may be encountered with the need to recompile the entire program every time a change is made by the Editor, and this could be time-consuming for large programs. It could however be avoided by a process of incremental compilation, as discussed by Atkinson et al.,<sup>151</sup> but for teaching purposes, when the programs are normally short, repeated recompilation is probably the better approach.

#### 2.6. FORMAL LANGUAGES IN CHEMISTRY AND INFORMATION WORK

Formal language theory has been applied in a number of areas in chemistry and information work: at the simplest level the interactive search languages used in online bibliographic retrieval systems<sup>152, 153</sup> have grammars which can be described by the methods developed by Chomsky. For the most part, they are Type 3 (Regular) languages, with trivial syntax analysis.

Some more sophisticated query languages have also been developed for specific applications, such as MQL (Medical Query Language)<sup>154</sup> which is designed to allow input of queries to a database in something approximating to natural language.

Specialised descriptive languages have been developed for use with chemical synthesis planning programs. In these programs the computer, upon being presented with a "target" chemical structure, is able by use of a database of chemical reactions

called **transforms** to suggest possible synthesis routes leading to the target.

Formal languages have been developed for the description of the **transforms**, two such being CHMTRN<sup>155</sup> used by the LHASA (Logic and Heuristics Applied to Synthetic Analysis)<sup>156, 157</sup> program, and ALCHEM (A Language for CHEMistry)<sup>158</sup> used by the SECS (Simulation and Evaluation of Chemical Synthesis)<sup>159, 160</sup> program, which is historically an offshoot of LHASA. Figure 2.2 illustrates the description of a **transform** using ALCHEM. Each **transform** contains information which enables the computer to decide whether or not it is applicable to the synthesis of a particular target molecule.

Both languages have been designed to represent the **transform** in a manner which remains reasonably intelligible to a chemist, yet is amenable to computer analysis, and "compiler" programs have been written for them. Both have a fairly strict line format, and their structure is more akin to that of Fortran than those of more modern languages such as Algol and its descendents; their grammars are not formalised by production rules or syntax diagrams.

Line notations used for the representation of chemical structures as strings of alphanumeric symbols can be regarded as formal languages, and some success has been achieved in writing a context-sensitive grammar for the Wiswesser notation.<sup>161</sup> Lin

```

1 TYPE PATTERN
2 ; PROXIMITY GUIDED EPOXIDATION
3 ; ALCOHOL GROUP CIS TO EPOXIDE ON RING
4 ; REF: E. COLVIN, J CHEM SOC PERKIN I 1989 (1973)
5 ; CHEM COMM 858 (1971), HOUSE P. 305
6 EPOX
7 O—C—C—@1<1, 3, 2>/
8 PRIORITY 0
9 CHARACTER ALTERS GROUP
10 ; CHECK IF STEREOCHEMISTRY IS IMPORTANT
11 IF STEREOCENTER IS CARBON OFFPATH THEN ; IT IS IMPORTANT
12 BEGIN IF ALCOHOL IS WITHIN GAMMA TO ATOM 2 (1) THEN
13 BEGIN IF BOND 1 AND (1) ARE CIS THEN ADD 50
14 ELSE KILL ;EPOXIDATION WOULD HAVE WRONG STEREOCHEM
15 IF (1) IS ONRING OF SIZE 5-6 THEN ADD 50
16 DONE
17 IF NITRILE IS EPSILON TO ATOM 2 (2) THEN
18 BEGIN IF BOND 1 AND (2) ARE TRANS THEN ADD 30
19 ELSE SUBT 30 ;EPOXIDE TRANS TO NITRILE IS FAVORED
20 DONE
21 DONE
22 CONDITIONS SLIGHTLY OXIDIZING
23 DELETE ATOM 1
24 MAKE BOND FROM ATOM 2 TO ATOM 3
25 END
26 COMPLETE

```

Figure 2.2: ALCHEM description of a chemical transform (from Wipke et al.<sup>160</sup>)

et al. have also written a compiler which performs automatic syntax<sup>162</sup> analysis on their Separate Feature Linear Notation (SEFLIN).<sup>163</sup>

Formal language theory has also been applied in chemistry outside the area of artificial language design. Fehder and Barnett<sup>164</sup> suggested in 1965 that the principles of syntactic analysis could be applied to the analysis of molecular formulae, providing a means for determining the validity (grammatical correctness) of a given molecular formula, and other authors have followed up this approach.<sup>165-167</sup> Similar applications have been made in the analysis of nomenclature.<sup>168-170</sup>

Rankin and Tauber<sup>171, 172</sup> have applied formal language theory to the full topological representation of chemical structures, developing generative grammars based on production rules for certain classes of molecule; such grammars are also discussed by Whitlock.<sup>173</sup>

In their second paper<sup>171</sup> Tauber and Rankin suggested that sets of grammar rules could be used for compact storage of groups of related structures, such as leucine esters, different rules being used for the generation of the constant and variable parts of the structure. A similar approach to the storage of generic structures was later taken by Krishnamurthy and Lynch.<sup>15, 16</sup>

More recent work by Welford<sup>67, 174</sup> has extended the range of structure types that can be generated and recognised by formal grammars, and has formed a cornerstone of the research on generic structure representation at Sheffield University<sup>174-177</sup> of which this Thesis describes a part.

## CHAPTER 3

### THE INPUT LANGUAGE

---

"My language is plain."

Bret Harte (1836-1902)

Chapter 1 has surveyed the various types of expression found in generic structure descriptions in patent specifications and abstracts, and has outlined the reasons for the development of a special input notation, or language, for the description of such structures which will be intelligible to a chemist, information scientist or patent agent, yet sufficiently formalised for automatic analysis by computer, using the principles discussed in Chapter 2.

The language described in this Chapter, GENSAL, may be used to represent a generic structure unambiguously (in order that an

unambiguous internal representation may be generated from it), and it has been designed to conform as far as possible to the type of description commonly found in chemical patent specifications. It is thus a formalised version of the generic structure description of patent specifications and abstracts: aspects of its formal grammar are described in Section 3.11., and as with many modern programming languages the grammar of GENSAL is expressed as a series of syntax diagrams, shown in Appendix 1. Throughout the text of this thesis, the syntactic metasymbols of GENSAL used as headings for the syntax diagrams are shown underlined.

### 3.1. GENERIC STRUCTURE DESCRIPTION USING GENSAL

The basic layout of generic structure descriptions in patent specifications and abstracts, as discussed in Section 1.2., is retained in GENSAL, one sentence of which describes one generic structure. Syntax Diagram 21 shows that the overall description of a structure has an introductory heading part, containing a reference number, and a **structure diagram** for the constant part of the structure which is followed by a series of statements, separated by semicolons; the sentence ends with a full stop. Figure 3.1 shows a simple generic structure and its GENSAL representation which, as can be seen, remains readily intelligible to a chemist.

The plethora of symbols used for structural and multiplicative

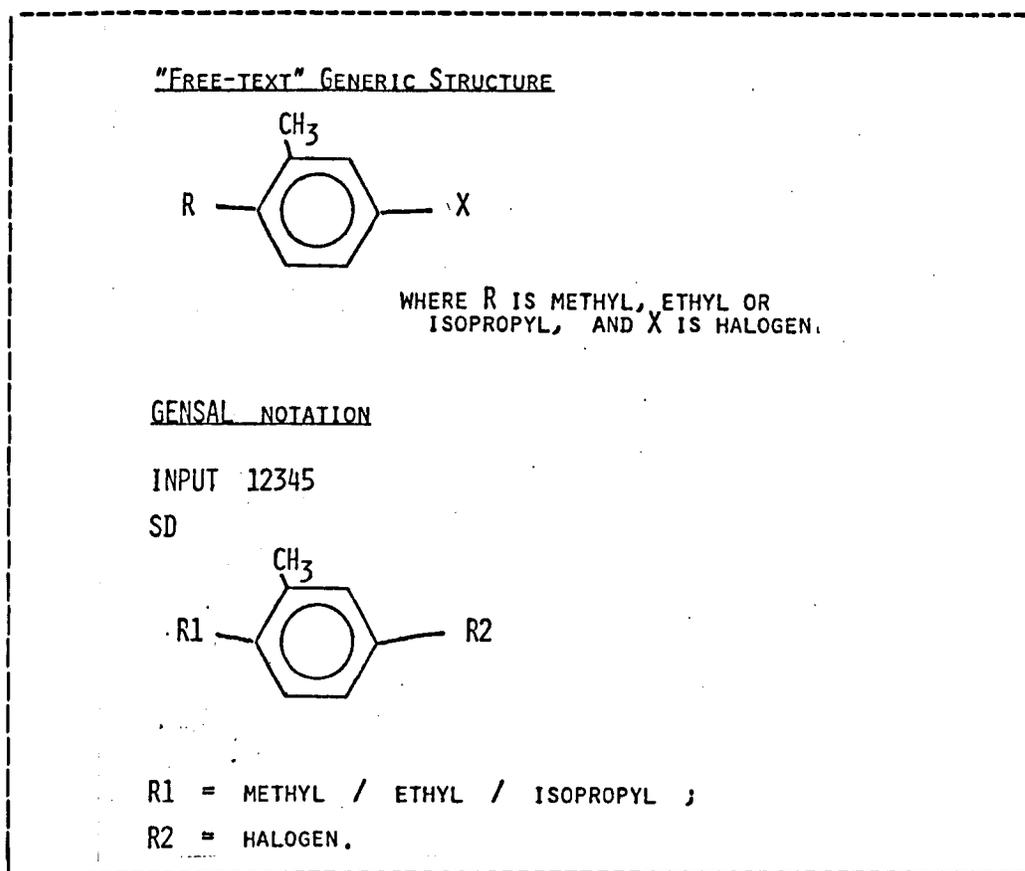


Figure 3.1

variables is reduced to two standard series: R1, R2, R3 etc. for structural variables (called substituents in GENSAL), and M1, M2, M3 etc. for multiplicative variables (called multipliers), as shown in Syntax Diagrams 3 and 4.

Variables in a GENSAL sentence must be introduced ("declared"), normally by appearing in a **structure diagram**, before being given values ("defined") in terms of chemical nature and position for substituents, and of selectors (giving integer ranges) for multipliers.

The definition of substituents and multipliers takes place in

assignment statements, which contain facilities for assigning the same set of alternatives to groups of substituents or multipliers simultaneously (with both independent and non-independent selection of the alternative values) or for assigning to substituent combinations (forming an extra ring). The substituent value may be given in several different ways, and there is scope for indicating the position at which the substituent is attached, and any further substitution on it, down to any level.

Conditional definitions are indicated by IF and RESTRICT statements. The former allow the use of one of two alternative subordinate statements according to whether a condition involving substituents and multipliers already defined is TRUE or FALSE. The latter impose such conditions on the alternatives given in earlier assignment statements, allowing only those combinations of alternatives that result in the condition being TRUE.

The next nine Sections of this Chapter give a detailed description of the language, allowing a full understanding of the GENSAL notations for the actual patent examples shown with the Derwent Abstracts of the original specifications in Figures 3.2 to 3.11. A comprehensive instruction manual for GENSAL, with further examples, has been prepared by Hill. <sup>178</sup>

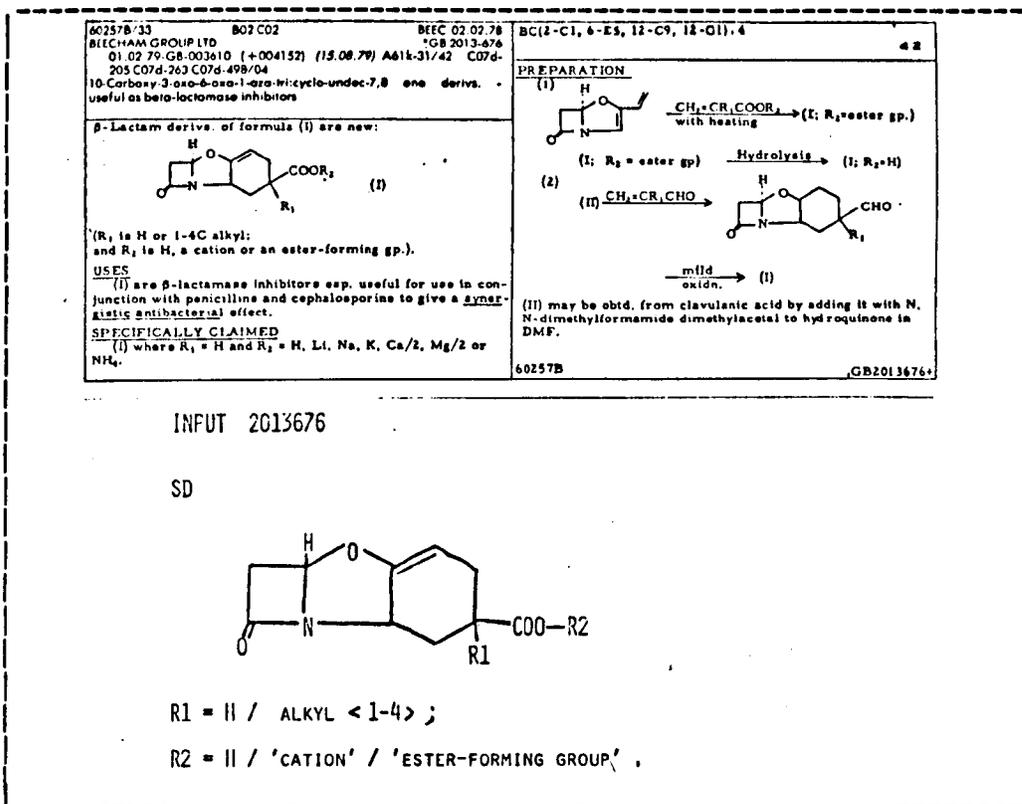


Figure 3.2

### 3.2. STRUCTURE DIAGRAM INPUT

As a whole, GENSAL is intended to be independent of any given computer system, and its high degree of readability makes it suitable as a means of describing generic structures manually, just as the programming language Algol is often used for the manual description of algorithms.

Nevertheless, certain aspects of GENSAL are intended to be implementation-dependent, and the most important of these are the structure diagrams which form an integral part of the language.

Any suitably-modified chemical structure graphics system might be used for their input, with a routine to convert its output into the connection table format used in the internal representation of the structure.

In the implementation described in Chapter 5 a modification of the program developed by Feldmann and others<sup>179</sup> and used in the Crystal Structure Search and Retrieval (CSSR)<sup>180</sup> and National Institutes of Health / Environmental Protection Agency (NIH/EPA)<sup>181</sup> substructure search systems, is being used. This is far from ideal, but has the advantages that it was provided free, and uses standard lineprinter characters in its display routines, and thus does not require any special hardware.

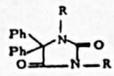
It is possible that an operational system might use a micro-computer as an intelligent terminal for the mainframe on which the bulk of the structure processing and searching would be carried out, and that the microcomputer would handle the chemical structure graphics locally, transmitting and receiving connection tables for each diagram.

A structure graphics system used with GENSAL requires certain features not found in all such systems. There must be a facility for defining nodes of the diagram as substituents (with the correct syntax) as well as as atoms of different types, and also a facility for applying multipliers (with the correct syntax) to nodes defined as substituents.

It must be possible to show that a particular node is connected back to a previously-defined part of the structure: in the modified Feldmann program used for the present work, this is achieved by attaching such an "apical" node to a dummy node, whose atomic type is given as "\*".

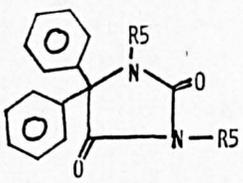
It must be possible to show that a particular node is attached to one of the other nodes in the diagram, without specifying which. In patent specifications, and general chemists' usage, this is usually achieved by the convention of a bond drawn into the centre of a ring, but in the modified Feldmann program it is done by attaching the variably-positioned node to a dummy node of atomic type "#", which indicates that it may be attached to any other node in the diagram with sufficient spare valencies. If it is desired to restrict the available positions to, for example, those in a particular ring, this should be done by using a GENSAL position set (Section 3.6 below) elsewhere in the structure description.

Only nodes defined as substituents may have dummy "variable-position" nodes attached to them.

<p>615028/33 803 INTE: 22.04.77                  INTERX RES CORP 'US 4163-058                  22.04.77-US-790087 (31.07.79) Adik-31/41 C07d-233/74                  5,5-Diphenylhydantoin derivs. - useful as anticonvulsants,                  anti-epileptics and antiarrhythmics</p>	<p>D(5-B1M, 7-D9, 12-D4, 12-F1), 4 110</p>
<p>Hydantoin derivs. of formula (I) and their acid-addn. and base salts, 1-4C alkyl halide quat. salts and N-oxides are new:</p> <div style="text-align: center;">  <p>(I)</p> </div> <p>(R is H or -CHR, -XR<sub>2</sub>;                  R<sub>1</sub> is H or 1-7C alkyl;                  X is O or S;                  R<sub>2</sub> is COR<sub>4</sub>; and                  R<sub>4</sub> is the acyl residue of a naturally occurring protein amino acid; provided that one R must be other than H).</p> <p><b>USE/ADVANTAGE</b>                  (I) are anticonvulsants, antiepileptics and antiarrhythmics. Compared with phenytoin and related cpds., (I) have greater solubility and enhanced stability without undesirable side effects. Dose is 200-800 mg. daily in man, 30-200 mg. daily in dogs etc.</p>	<p><b>WIDER DISCLOSURE</b>                  (I) are stated to be new when (a) R<sub>1</sub> is also CCl<sub>3</sub>, CBr<sub>3</sub>, Cl<sub>3</sub>, Ph, Me<sub>3</sub>NCH<sub>2</sub>, CHO, PhOCH<sub>2</sub>, PhCH=CH, 2-furyl, pyridyl or R<sub>3</sub>-phenyl; R<sub>3</sub> is OH, Cl, Br, I, OMe, COOMe, NO<sub>2</sub> or OCOMe; (b) X is also NR<sub>1</sub>; (c) R<sub>2</sub> is also PO(OH)<sub>2</sub>; and (d) R<sub>4</sub> is also R<sub>3</sub>C<sub>6</sub>H<sub>4</sub>, pyridyl, pyridylmethyl, pyridyl N-oxide, the residue of any N-substit. amino acid (the substituent being a protecting gp. removable by hydrogenolysis or hydrolysis), the residue of an N,N-di-(1-5C)alkyl- or 4-7C cycloalkyl-amino acid, (CH<sub>2</sub>)<sub>n</sub>COOH, CH<sub>2</sub>OCH<sub>2</sub>COOH, (CH<sub>2</sub>)<sub>n</sub>COOMe, (CH<sub>2</sub>)<sub>n</sub>COOEt, (CH<sub>2</sub>)<sub>n</sub>CONR<sub>4</sub>, imidazolyl, 1-8C alkoxy, PhCH<sub>2</sub>O, PhO or O(CH<sub>2</sub>)<sub>n</sub>NR<sub>4</sub>; where n is 1-5 and R<sub>3</sub> and R<sub>4</sub> are 1-5C alkyl or NR<sub>4</sub>; R<sub>4</sub> is a heterocyclic gp.</p> <p><b>SPECIFICALLY CLAIMED</b>                  3-[N,N-Dimethylglycylloxymethyl]diphenylhydantoin (II) and its methanesulphonate and salicylate.</p> <p><b>PREPARATION</b>                  5,5-Diphenylhydantoin (III) is either first treated with R<sub>1</sub>CHX and then acylated with a reactive deriv. of R<sub>4</sub>COOH, or it is converted to a Na, K, Ca or Mg deriv. and this is</p>

INPUT 4163058

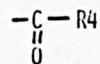
SD



R5 = H / SD  
 -CH-R3-R2  
 |  
 R1

R1 = H / ALKYL <1-7> ;

R2 = SD



R3 = SD -O- / SD -S- ;

R4 = 'ACYL RESIDUE OF NATURALLY OCCURRING PROTEIN AMINO ACID' ;

RESTRICT <1> R5 <> H .

Figure 3.3

3.3. SIMPLE ASSIGNMENT STATEMENTS

Assignment statements express the definition of a substituent or multiplier. For both types of variable, the substituents or multipliers being defined are shown on the left-hand side of an assignment operator (normally "="), and the possible values on

the right-hand side.

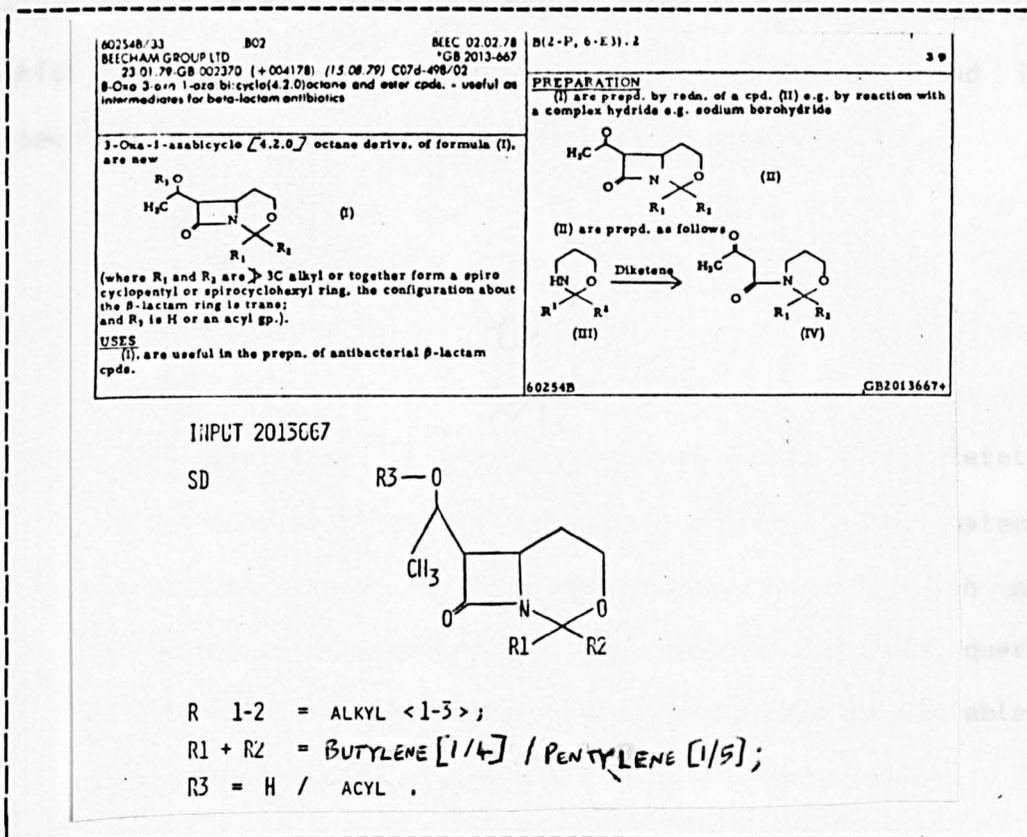
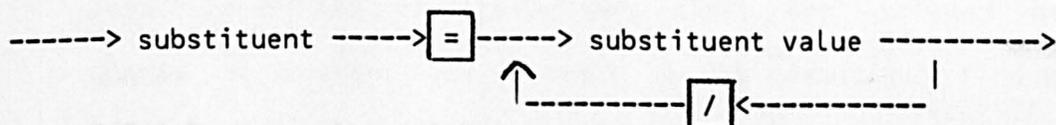


Figure 3.4

### 3.3.1. Substituent Assignments

The simplest form of assignment statement for substituents can be represented as follows:



(This is a simplified version of the relevant syntax diagrams.)

It allows a single substituent on the left-hand side to be

defined as having one of the values separated by the "/" delimiters on the right-hand side. Each alternative value may be given in one of five different forms, shown in Syntax Diagram 10, which correspond to the different types of expression found in specifications and abstracts and discussed in Section 1.2.2.

#### 3.3.1.1. Unknown Value

A "?" represents a substituent whose nature is completely unknown. This situation usually occurs with patent expressions such as "optionally substituted", with no indication of the nature of the substitution. In query structures it might also be used as a value for variables indicating the unspecified parts around a substructure.

#### 3.3.1.2. Structure Diagram

The substituent is defined by a structure diagram, which is input in exactly the same way as the main structure diagram for the constant part of the structure, and may, of course, have further substituents declared within it.

#### 3.3.1.3. Nomenclatural Terms and Expressions

**Specific nomenclatural terms** represent a single chemical

entity, and include terms such as "chloro", "methyl", "pyridyl" and "cyclohexyl". Simple linear formulae such as "CN", "COOH", and "NH<sub>2</sub>" are also regarded as **specific nomenclatural terms**.

Essentially this is a shorthand method of inputting a **structure diagram**: an operational system might have sophisticated routines for nomenclature translation and linear formula analysis, though development of these has not formed part of the present work. At a simpler level, when a GENSAL sentence is being interpreted by computer, a file may be searched for a record of the structure of, e.g. phenyl, and if no entry is found a suitable message be printed at the terminal. Such a file may also be able to simplify compound terms such as "halogen" and "alkali metal". This is the approach used in the current implementation.

**Homologous series terms** describe classes of structural entities, and the **parameter list** which follows the term may impose restrictions on the variety of structures covered. **Parameter lists** are discussed more fully in Section 3.5.

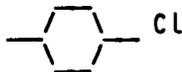
Verbal expressions that do not correspond to any specific structure or structurally-defined class are enclosed in quotes to prevent any attempt by the computer to find a structure record.

The following are simple examples of assignment statements:

R1 = methyl / ethyl / cyclohexyl ;

R2 = 'electron-withdrawing group' ;

R3 = SD



/ cyclohexyl ;

R4 = ?

and further examples may be found in the GENSAL notations for patent examples shown in Figures 3.2 and 3.3.

### 3.3.2. Multiplier Assignments

Simple multiplier assignments are of the form:

-----> multiplier -----> [=] -----> selector ----->

and enable multipliers to be assigned a range of integer values.

Such a range is defined using the integer range given in Syntax Diagram 2 enclosed in angle brackets. It may consist of a single integer, or a group of single integers or "range fragments" separated by commas. Each range fragment consists of two integers separated by a hyphen, and represents all the integers from the lower to the upper inclusive. The last integer before the end of the selector (immediately before the right angle bracket) may be

followed by a hyphen without a second integer, in which case all the integers from the bound upwards are included.

Thus the selector <0-6,8,12,15-19,23-25,31,43-> includes the integers 0, 1, 2, 3, 4, 5, 6, 8, 12, 15, 16, 17, 18, 19, 23, 24, 25, 31, 43, 44, 45, 46, etc., potentially up to infinity.

As can be seen from Syntax Diagram 1, negative integers are not allowed in GENSAL, and there is also a requirement that the values in an integer range (Syntax Diagram 2) must increase from left to right.

An example of a simple multiplier assignment is

$$M1 = \langle 3-5 \rangle$$

and further examples may be found in Figures 3.5 and 3.10.

#### 3.4. MORE COMPLEX ASSIGNMENTS

The full syntax diagram for assignment statement (No. 16), and those with which it is defined, allow much more complicated assignments to be concisely represented.



substituent values describe only the atoms that are added to the structure, i.e.:



The former has the advantage of being more consistent with normal patent usage, but the latter is simpler to implement, and has been chosen for the present work.

Only two substituents may be combined in this way, and both must be singly connected in their independent existence.

Figure 3.4 shows an example of substituent combination in a patent; in this case the possibility of combining R1 and R2 is alternative to their being separate singly-connected radicals.

#### 3.4.2. Group Assignment Statements

As a convenient "short-cut" several substituents or multipliers can be defined simultaneously:



i.e. R1 and R2 are both defined by the three alternatives shown. (Note that in this case the three nomenclatural terms are being used to represent singly-connected radicals, substituted on the constant part of the structure, in contrast to their use above

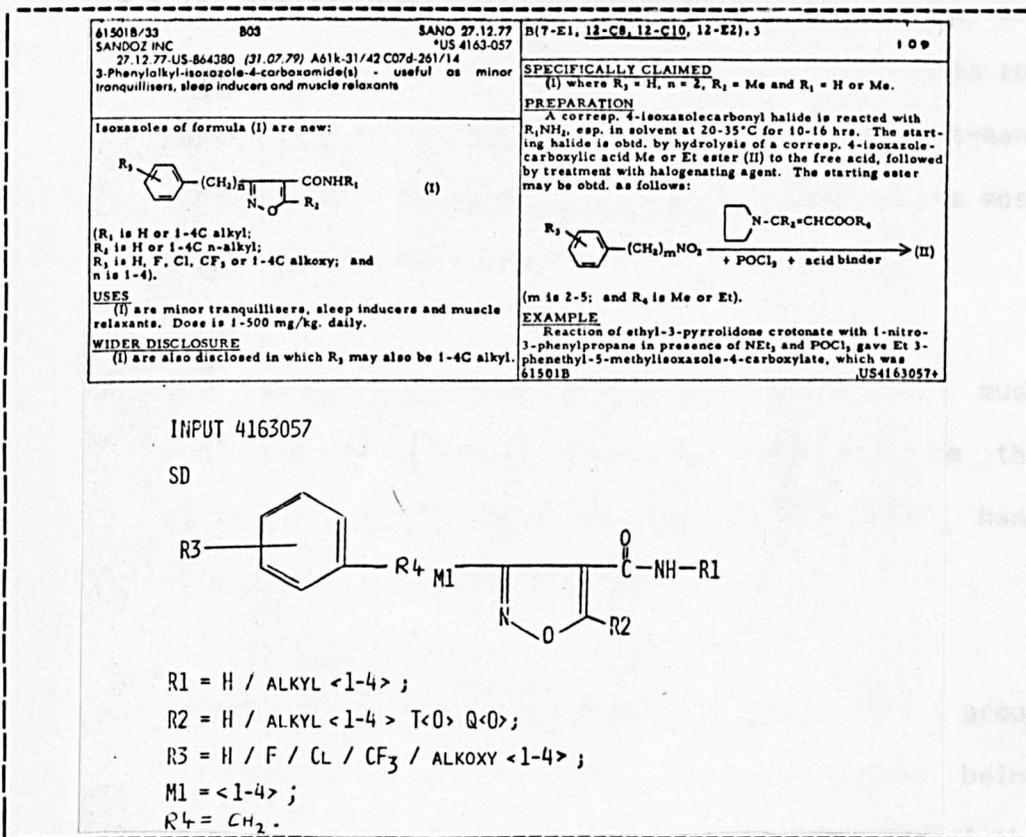


Figure 3.6

with a substituent combination.) The integers used to list the substituents being defined are arranged in the syntax for an integer range.

### 3.4.2.1. Assignment Operators

The five available assignment operators shown in Syntax Diagram 16 have different meanings, which may be useful when several substituents or multipliers are being defined together in a group assignment, and their values are not necessarily independent.

= : The substituents or multipliers in the group are independently selected from the alternatives in the substituent definition or selector on the right-hand side of the assignment statement. This is the most commonly-encountered operator.

S= : All the substituents or multipliers in the group must have the same value, which is selected from the substituent definition or selector on the right hand side of the assignment statement.

D= : Each of the substituents or multipliers in the group must have a different value, all the values being selected from those on the right-hand side of the statement.

\$= : Not all the substituents or multipliers in the group may be the same (which, using "=", they could be), but they need not all be different.

#= : Not all the substituents or multipliers in the group may be different (i.e. at least two must be the same), but they need not all be the same.

Examples of such simultaneous assignments are as follows:

a) R1-3 = phenyl/ cyclohexyl / cyclopentyl ;

R1, R2 and R3 can be independently either phenyl, cyclohexyl or cyclopentyl. There are thus 27 (3 x 3 x 3) possible permutations, assuming that there are no symmetry considerations involved.

b) R4-6 D= phenyl / cyclohexyl / cyclopentyl;

R4, R5 and R6 must be different, each being selected from the possibilities phenyl, cyclohexyl and cyclopentyl. There are thus 6 (3 x 2 x 1) possible permutations.

c) R7,9-10 S= phenyl / cyclohexyl /cyclopentyl;

R7, R9 and R10 must be the same. There are only three possible permutations (all phenyl, all cyclohexyl or all cyclopentyl).

d) R11,13,15 \$= phenyl /cyclohexyl / cyclopentyl;

R11, R13 and R15 are not all the same, each being otherwise selected from the possibilities given. This leaves 24 possible permutations, there being three ways in which all can be the same.

e) R16-18 #= phenyl / cyclohexyl / cyclopentyl;

R16, R17 and R18 are not all different, but are

otherwise selected from the available possibilities. Here there are 21 possible permutations, as there are six ways in which the three may be all different.

#### 3.4.2.2. Selected Group Assignments

A group assignment statement can begin with a selector which allows just some of the substituents or multipliers in the group to be assigned values from the substituent definition or selector on the right-hand side of the assignment statement. e.g.:

```
<2-3> R1-5 = phenyl / cyclohexyl / cyclopentyl;
```

means that 2 or 3 of the group of substituents R1, R2, R3, R4 and R5 are independently selected from the list phenyl, cyclohexyl and cyclopentyl, the others remaining undefined at this point in the GENSAL structure description.

### 3.5. HOMOLOGOUS SERIES IDENTIFIERS AND GRAMMARS

Certain terms used in generic structures cover a range of specific substructures all of which are alternative to each other at that point. The most common example is the term "alkyl" which covers all rooted acyclic substructures containing carbon and hydrogen only, without any saturations.

Welford<sup>67, 174</sup> has developed chemical grammars to deal with this type of expression, and not only are they applicable to terms such as alkyl and alkenyl, which are commonly understood by chemists as "homologous series terms", but they may also be applied to many less precise terms which are none the less "structurally recognisable" -- that is, terms which encompass a range of substructures that have a particular structural feature in common, such as "aryl" and "heterocyclic". Each valid **homologous series identifier** is associated with a list of parameters to the chemical grammars, the values of each parameter being defined by means of a selector.

As can be seen from the syntax diagram for parameter (No. 7), the parameter may be indicated by a Parameter Identifier or a substituent enclosed in quotes. The standard Parameter Identifiers cover features such as atom count, branch points etc., and are shown with their meanings in Table 3.1, though it is possible that the list may be modified as the chemical grammars are further developed. Non-standard parameters, shown by substituents in quotes, cover such features as interruptions in the chain, and substitutions on it. Thus for the **homologous series identifier** "alkyl", all the parameters will be zero, except for C, T, Q and P, which can take on any (mutually consistent) values.

Syntax Diagram 10 allows a parameter list to follow a **homologous series identifier**, thus more closely defining any of the standard parameters or introducing non-standard ones. This results in

C	Carbon count
E	double bonds (alkEne)
Y	triple bonds (alkYne)
Q	Quaternary branches
T	Ternary branches
RC	number of Rings
RN	Number of Ring atoms
RS	number of Ring Substitutions
RF	number of Ring Fusions
RA	number of Aromatic Rings
Z	number of heteroatoms

Table 3.1: Parameter Identifiers and their meanings

expressions like:

alkyl C<3-8> T<1-2>

which indicates alkyl groups containing between 3 and 8 carbon atoms, with 1 or 2 ternary branching atoms. The parameters may appear in any order, or be absent altogether, in which case their default values will be the widest possible range compatible with those parameters that are present (including any implicit in the **homologous series identifier** itself). In view of the way in which the Parameter Identifier "C" (for carbon count) occurs almost every time a **homologous series identifier** is used a shorthand has been introduced whereby the C may be omitted, provided this is the first parameter in the list. e.g.

alkyl <1-4>

alkyl <3-8> T<1-2>

Further examples of **homologous series terms** and parameter lists are:

a) cycloalkyl <10-20> RC<2->

(Between 10 and 20 carbon atoms, and at least two rings.)

b) alkyl <3-12> 'R5' <0-1>

(Between 3 and 12 carbon atoms, and 0 or 1 occurrences of R5.)

c) carbacyclic <6-10> E <1->

(Between 6 and 10 carbon atoms, and at least one double bond (number of triple bonds not specified). The term "carbacyclic" is used to indicate acyclic hydrocarbons.)

The assignment statements for R2 and R3 in Figure 3.5 include **homologous series identifiers** and parameter lists; in the case of R2 the specification of no ternary or quaternary branching atoms is equivalent to the statement in the Derwent Abstract that it is an n-alkyl group.

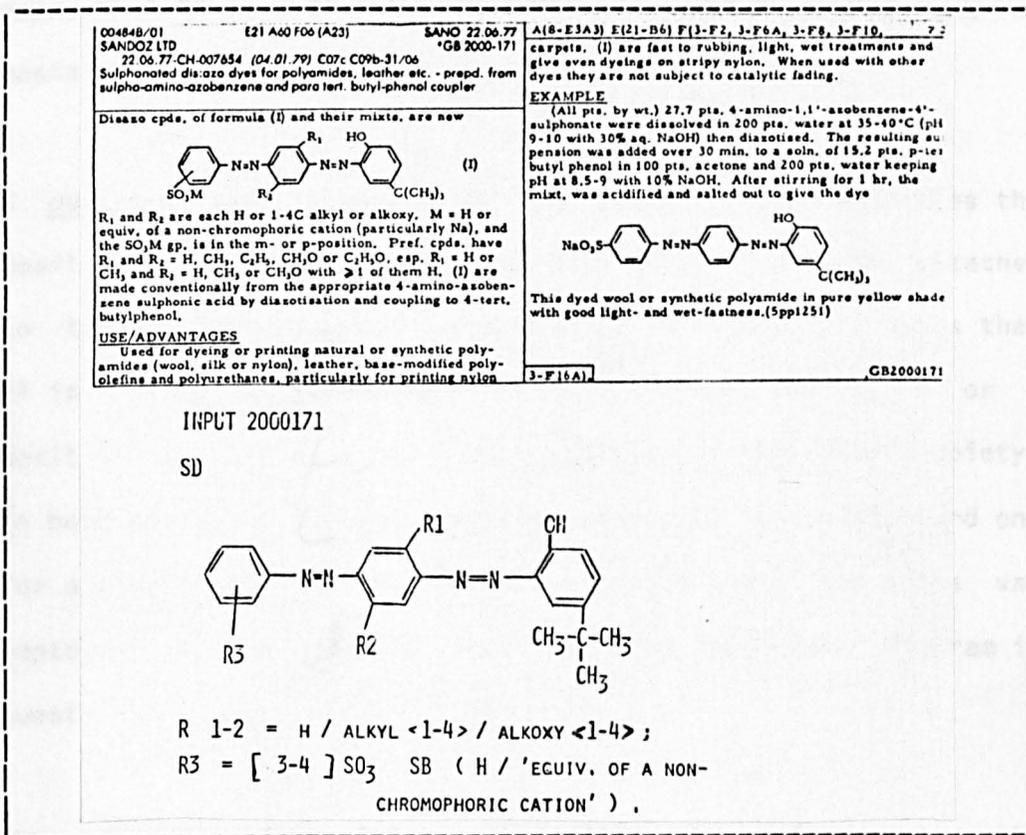
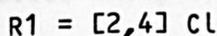


Figure 3.7

### 3.6. POSITION SETS

Syntax Diagram 14, for definition element (of which substituent value is a simple case), allows the inclusion of some information about the position(s) of substituents being defined. A position set at the beginning of a definition element indicates the position(s) in the constant structure at which the substituent(s) currently being defined may be attached. Thus



means that R1 is a Cl group attached in either position 2 or position 4 of the constant structure.

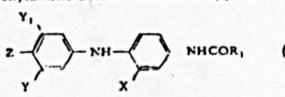
A position set following a substituent value indicates the position(s) in the substituent through which it may be attached to the constant structure. The example in Figure 3.12 shows that R1 is a nicotinic acid moiety attached through its 2, 4 or 6 position to the 2, 3, 4, 5 or 6 position of the toluene moiety. In both cases the numbering system referred to is the standard one for a nomenclatural term, or whatever numbering of the atoms was employed in the graphic input of the **structure diagram** in question.

The appearance of position sets after the reference to the structure to which they apply allows immediate automatic checking on the availability of the specified positions in the structure in question, which can be valuable in the machine processing of GENSAL.

Figures 3.6 and 3.7 show patent examples involving position sets. The extent to which position sets need to be used in GENSAL notations may depend on the facilities available in the graphics system being used for indicating alternative positions of attachment (especially to rings).

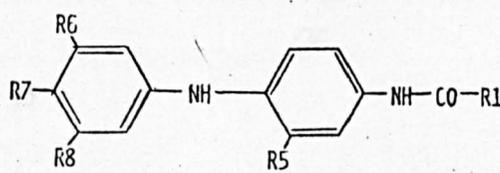
For doubly-connected substituents, or combinations of singly-connected ones, it may be necessary to specify pairs of positions in the position set. In this case, the positions in each pair are

separated by a "/", and the alternative pairs by commas. The order of the positions in such pairs is significant, and this may be important if the bond orders of the two connections are different, or if the substituent value is not symmetrical.

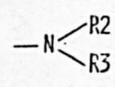
<p>3MILLION CO 22 01 79 US-005633 (+761515) (01.01.80) A01n-09/12 C07-103/73 4-Acylamino di-phenylamine derivs. - useful as herbicides, esp. for post/emergent application</p>	<p>US 4181-519 R<sub>1</sub> = H, 1-4C alkyl, F, Cl or Br; provided (a) when Z = H, X is not haloalkyl; (b) when R<sub>1</sub> is 1-R<sub>4</sub>-cyclopropyl, at least 1 of Y, Y<sub>1</sub> and Z is other than H; (c) when R<sub>1</sub> is alkyl at least 2 of Y, Y<sub>1</sub> and Z are other than H).</p>
<p>Intermediate Priority: 10.2.78 - US-876593. Diphenylamine deriva. of formula (I) are new:</p>  <p>(X is F, Cl, Br or 1-6C alkyl, opt. substd. by ≥ 1 halo; Y and Y<sub>1</sub> = H, F, Cl, Br or alkyl or alkoxy opt. substd. by ≥ 1 halo; Z is H, F, Cl, Br, 1-6C alkyl, alkoxy, alkylthio, alkylsulphinyl or alkylsulphonyl (all opt. substd. by ≥ 1 halo) or is R<sub>2</sub>R<sub>3</sub>N; R<sub>2</sub> and R<sub>3</sub> = H or ≥ 6C alkyl or cycloalkyl; R<sub>2</sub> may also be 1-6C alkoxy; R<sub>4</sub> = 1-4C alkyl or 1-R<sub>4</sub>-cyclopropyl;</p>	<p>USE (I) are herbicides (some selective), esp. for post-emergent use. The usual application rate is 0.1-10 lbs./acre in standard formulations. PREPARATION By acylating the appropriate 4-aminodiphenylamine (II) with R<sub>1</sub>COHal in presence of usual solvents and acid acceptors. (II) are made by condensing a Y, Y<sub>1</sub>, Z-substd. aniline with a 3-X-4-chloronitrobenzene, then reducing the resulting 4-nitrodiphenylamine. EXAMPLE 17.8g. 4-amino-2,3'-bis(trifluoromethyl)-4'-chlorodiphenylamine and 5g. Et<sub>3</sub>N were dissolved in 100 ml. THF, then treated dropwise with 6g. 1-methylcyclopropanecarboxonyl chloride. The mixt. was refluxed for 1 hr., poured into ice-water and extracted with 3 x 150 ml. ether. The ex-</p>

INFUT 4181519

SD



R1 = ALKYL <1-4> / CYCLOPROPYL SB [1] R4 ;  
R4 = H / ALKYL <1-4> / F / CL / BR ;  
R5 = F / CL / BR / ALKYL <1-6> SB <0-> HALO ;  
R 6,8 = H / F / CL / BR / ( ALKYL / ALKOXY )  
SB <0-> HALO ;  
R7 = H / F / CL / BR / ( ALKYL <1-6> / ALKOXY <1-6> /  
ALKYLTHIO <1-6> / ALKYL SULPHINYL <1-6> /  
ALKYL SULPHONYL <1-6> ) SB <0-> HALO / SD



R 2-3 = H / ALKYL <-6> / CYCLOALKYL <-6>;  
R3 = ALKOXY <1-6>;  
IF R1 = CYCLOPROPYL SB [1] R4 THEN  
RESTRICT <1-> R 6-8 <> H ;  
IF R1 = ALKYL THEN RESTRICT <2-> R 6-8 <> H ;  
IF R7 = H THEN RESTRICT R5 <> ALKYL <1-6> SB <1-> HALO .

Figure 3.8

### 3.7. NESTED SUBSTITUTION

As was stated in Chapter 1, the variable parts in a generic structure may themselves be further substituted to any level. GENSAL is able to show this clearly and concisely by means of the mutually recursive Syntax Diagrams 14 and 15 in which parentheses are used in expressions involving the four substitution operators "/", "&", "SB" and "OSB" to remove any possible ambiguity.

The operators respectively represent "or", "and", "substituted by" and "optionally substituted by", and are evaluated in the order "&", followed by "SB" and "OSB" (ranking equally and evaluated from left to right), followed by "/". Expressions in parentheses are evaluated before "&".

This precedence order has been adopted because it appears to provide the most natural form for complex expressions: AND is conventionally evaluated before OR in Boolean expressions, and the intermediate positioning of SB and OSB allows several substitutions to be made on each alternative without the use of parentheses.

The following expression includes examples of the use of all four operators:

```
R1 = phenyl sb (Cl / Br / I) & nitro /  
      (cyclohexyl / cyclopentyl)  
      sb (amino / pyridyl osb methyl & methoxy) /  
      naphthyl
```

The expression indicates that R1 has three possible basic alternatives:

1. A phenyl group substituted both by
  - a) either Cl or Br or I
  - and by b) nitro
2. A cyclohexyl or cyclopentyl group substituted by either
  - a) an amino group
  - or by b) a pyridyl group, itself optionally further substituted both by methyl and by methoxy (i.e by both or by neither)
3. A naphthyl group, not further substituted.

The examples in Figures 3.8 and 3.9 include assignment statements involving parentheses to indicate further substitution, and that in Figure 3.9 also shows how it may be necessary to alter the way in which the generic structure is expressed in the original specification or abstract in order to encode it in GENSAL. The design of GENSAL is intended to minimise the need for such alterations, but occasionally they do become necessary.

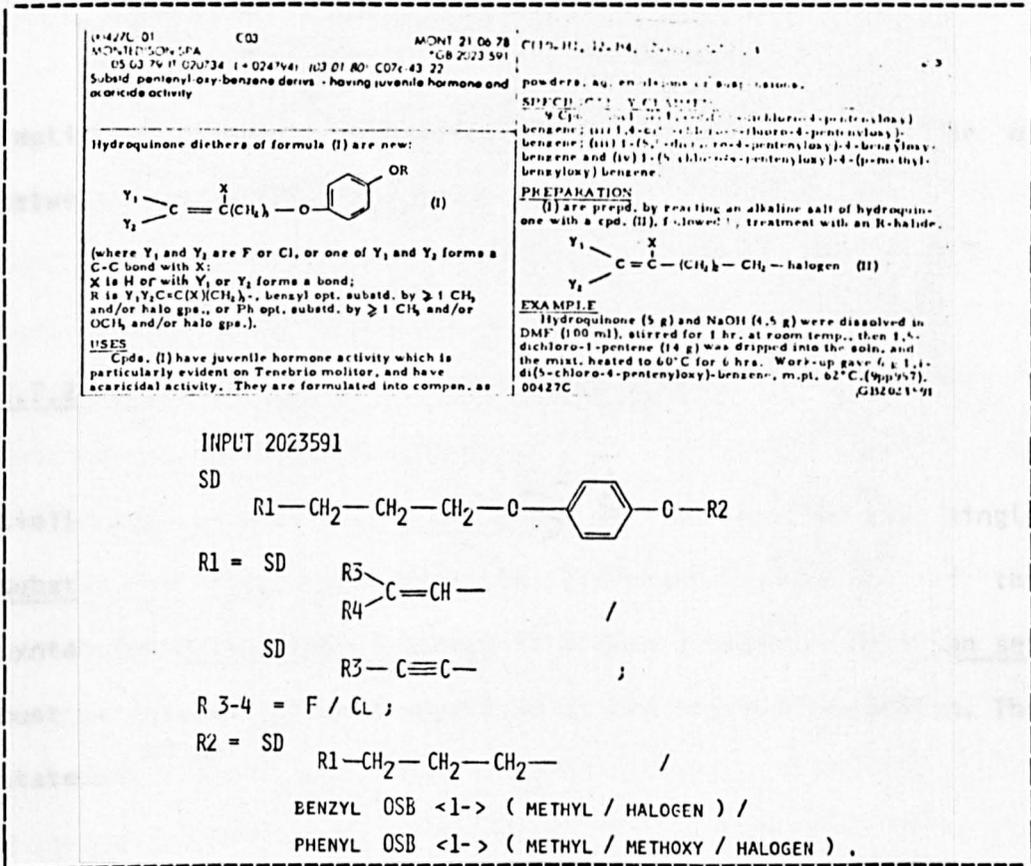


Figure 3.9

3.7.1. Selectors in Definition Expressions

A selector at the start of a definition element indicates multiple occurrences of whatever follows. Since the main part of the definition element may be a parenthesised definition expression, it is possible to have several such selectors applying to the same substituent value: in this case their effects are multiplied, with the result that the expression

<2-4> ( <2> methyl / ethyl )

implies the presence of between 4 and 8 methyl groups, or of between 2 and 4 ethyl groups.

### 3.7.2. Position Sets in Definition Expressions

Similarly, several position sets may be applied to a single substituent value, appearing in different recursions of the syntax for definition elements. Here each successive position set must be a subset of that specified at the previous recursion. The statement

R1 = [2-6] (methyl / [3,5] hydroxy)

is therefore valid, whereas

R1 = [2,4,6] (methyl / [3,5] ethyl)

is erroneous.

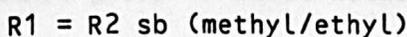
### 3.7.3. Substituents as Substituent Values

It may be convenient, and is occasionally found in specifications and abstracts, to define one structural variable in terms of

<p>059418/03 OLIN CORP 20.05.77-US-799970 (02.01.79) C10m-03/46 hydraulic and heat-transfer fluids - contg. alkoxy-siloxane(s)</p>	<p>E11 H08 OLIN 20.05.77 *US 4132-664</p>	<p>E(5-E) H(8-DS). S e o (I) and (II) have good hydrolytic stability, high lubricity and low viscosity indices. <b>DETAILS</b> The prepn. of cpds. (I) and (II) is described in US 392429. The pref. cpds. are those where a = 2; M is a 4-12C hydrocarbon radical; R is H, 1-8C alkyl or alkenyl, or 6-14C aryl or aralkyl; R' is hindered 4-12C alkyl. The fluids can comprise (I) and/or (II) alone or diluted, e.g., with glycols or glycol mono- or diethers. They can also contain conventional additives. <b>EXAMPLE</b> (II; a = 2, M = CH<sub>2</sub>CH<sub>2</sub>, R = Me, R' = n-Bu) has a pour point of &lt; -40°F, a VI of 339, a 4-ball wear scar of 0.66 mm and a flash point of 420°F.(7pp)67).</p>
<p>The use of hydraulic and heat-transfer fluids contg. alkoxy-siloxanes of formula (I) and/or (II)</p> <div style="display: flex; justify-content: space-around;"> <div style="text-align: center;"> <math display="block">M \begin{bmatrix} \text{Si}(\text{OR}^1)_2 \\   \\ \text{O} \\   \\ -\text{O}-\text{Si}-\text{O}-\text{Si}(\text{OR}^1)_2 \\   \\ \text{O} \\   \\ \text{Si}(\text{OR}^1)_2 \\ \text{(I)} \end{bmatrix}_a</math> </div> <div style="text-align: center;"> <math display="block">M \begin{bmatrix} \text{Si}(\text{OR}^1)_2 \\   \\ \text{O} \\   \\ -\text{O}-\text{Si}-\text{R} \\   \\ \text{O} \\   \\ \text{Si}(\text{OR}^1)_2 \\ \text{(II)} \end{bmatrix}_a</math> </div> </div> <p>where a = 2-4; M is an opt. substd. branched or straight-chain hydrocarbon radical; and R and R' are each independently selected from H, alkyl, alkenyl, aryl and aralkyl, provided that at least the majority of R' are sterically hindered <math>\geq 3\text{C}</math> alkyl) is claimed.</p> <p><b>ADVANTAGES</b></p>		<p>US4132664</p>
<p>INPUT 4132664</p> <p>SD</p> <p style="margin-left: 40px;">R1 — (R2)M1</p> <p>R1 = CARBACYCLIC SB ? ;</p> <p>R2 S= R3 / R4 ;</p> <p>R3 = SD</p> <div style="margin-left: 100px;"> <math display="block">-\text{O}-\text{Si} \begin{cases} \text{O}-\text{Si}-(\text{R5})_{\text{M2}} \\ \text{O}-\text{Si}-(\text{R5})_{\text{M2}} \\ \text{O}-\text{Si}-(\text{R5})_{\text{M2}} \end{cases}</math> </div> <p>R4 = SD</p> <div style="margin-left: 100px;"> <math display="block">-\text{O}-\text{Si} \begin{cases} \text{O}-\text{Si}-(\text{R5})_{\text{M2}} \\ \text{R6} \\ \text{O}-\text{Si}-(\text{R5})_{\text{M2}} \end{cases}</math> </div> <p>R 5-6 = H / ALKYL / ALKENYL / ARYL / ARALKYL ;</p> <p>M1 = &lt; 2-4&gt;;</p> <p>M2 = &lt; 3&gt; ;</p> <p>IF R2 = R3 THEN IF M1 = &lt; 2&gt; THEN RESTRICT &lt; 10-&gt; R5 = ALKYL &lt; 3-&gt;</p> <p style="margin-left: 80px;">ELSE IF M1 = &lt; 3&gt;</p> <p style="margin-left: 120px;">THEN RESTRICT &lt; 14-&gt; R5 = ALKYL &lt; 3-&gt;</p> <p style="margin-left: 120px;">ELSE RESTRICT &lt; 19-&gt; R5 = ALKYL &lt; 3-&gt;</p> <p style="margin-left: 40px;">ELSE IF M1 = &lt; 2&gt; THEN RESTRICT &lt; 7-&gt; R5 = ALKYL &lt; 3-&gt;</p> <p style="margin-left: 80px;">ELSE IF M1 = &lt; 3&gt;</p> <p style="margin-left: 120px;">THEN RESTRICT &lt; 10-&gt; R5 = ALKYL &lt; 3-&gt;</p> <p style="margin-left: 120px;">ELSE RESTRICT &lt; 13-&gt; R5 = ALKYL &lt; 3-&gt; .</p>		

Figure 3.10

understood that the substituent is on the highest level of substitution within the parenthesized expression. For example, in another. This is usually done where there is further substitution involved, as in



or

R1 = (phenyl/naphthyl) sb R2

Syntax Diagram 14 permits a substituent in place of a substituent value, and the substituent given may or may not already have been defined; the definition eventually given is treated as a parenthesised expression.

GENSAL additionally permits a substituent to be defined in terms of itself, as in

R1 = methyl sb R1

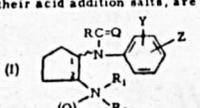
and it can be seen that this corresponds to an infinite-length polymer.

#### 3.7.4. Further Substitution on Parenthesised Expressions

Where one of the further substitution operators "SB" and "OSB" appears immediately after a parenthesised expression, it is understood that the substitution is made on the highest level of substitution within the parenthesised expression. For example, in the expression

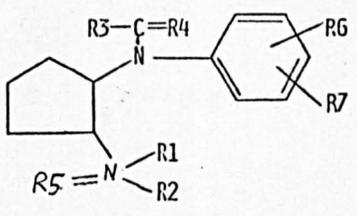
(phenyl / cyclohexyl sb methyl) sb Cl

the phenyl group is substituted by chlorine, and the cyclohexyl group by both methyl and chlorine; the methyl group is not further substituted.

<p>5744/B/78 B05 UPJO 30.11.76                  UPJOHN CO US 4159.340                  15 05 78 US 906429 (+ 746191) (26.06.79) CD7c-103/34                  Amino-cyclophatic amide(s) - possess CNS antidepressant                  properties and have a better therapeutic ratio than imipramine</p>	<p>B(10-A), 10-B2F, 12-C6, 3 106</p> <p>alkoxy and when Y is CF<sub>3</sub>, Z is H, when Y is 1 or 2C alkoxy                  and Z is H, the 1 or 2C alkoxy is in the 3-posn., when Y and                  Z are both halo or 1 or 2C alkoxy, they are present in the                  3- and 4-, or 3- and 5-posns).</p>
<p>Intermediate priorities: 15.3.77; 9.2.78-US-777599;                  876349.</p> <p>N-(2-Aminocyclopentyl)-N-acylanilides of formula (I) and                  their acid addition salts, are new</p>	<p>USE                  (I) possess potent CNS antidepressant properties and have                  a better therapeutic ratio than imipramine and long-acting                  activity which allows longer durations between administrat-                  ions, e.g. once a day. Dose is 1-100 mg/day.</p>
<p>(I)</p>  <p>(the wavy line indicates cis                  or trans configuration of                  the substituents in the 1 or 2                  posns. of the cyclopentyl                  ring;</p>	<p><b>SPECIFICALLY CLAIMED</b>                  3,4-Dichloro-N-2-(N-allyl-N-methylamino)cyclopentyl-                  propionanilide, 3,4-dichloro-N-2-(N-β-phenethyl-N-                  methylamino)cyclopentyl/propionanilide; 3,4-dichloro-N-2-                  (N-benzyl-N-methylamino)cyclopentyl/propionanilide; and                  their acid addn. salts.</p>
<p>p is 0 or 1;                  Q is O or S;                  R is 1-3C alkyl, vinyl, 3-6C cycloalkyl, ethoxy or methoxy-                  methyl;                  R<sub>1</sub> is H or 1-3C alkyl;                  R<sub>2</sub> is CH<sub>2</sub>C<sub>6</sub>H<sub>5</sub>, (CH<sub>2</sub>)<sub>3</sub>C<sub>6</sub>H<sub>5</sub>, 3-6C (allylic) alkenyl;                  Y and Z are each H, F, Cl, Br, CF<sub>3</sub>, 1 or 2C alkyl, 1 or 2C</p>	<p><b>PREPARATION</b>                  Cpd. (I), Q is O and p is 0) are prepd. by heating a cpd.                  of formula (II) and an anhydride of the appropriate carboxy-                  lic acid to form the N-acylated prod.</p>
	52472B US4159340

INPUT 4159340

SD



R1 = H / ALKYL <1-3>;  
 R2 = SD  
     - CH<sub>2</sub> -  /  
 SD  
     - CH<sub>2</sub> - CH<sub>2</sub> -  /  
 ALKENYL <3-6>;  
 R3 = ALKYL <1-3> / VINYL / CYCLOALKYL <3-6> /  
 ETHOXY / METHOXYMETHYL ;  
 R4 = O / S ;  
 R5 = <O-1> O ;

R 6-7 = H / F / Cl / Br / CF<sub>3</sub> / ALKYL <1-2> / ALKOXY <1-2> ;  
 IF R6 = CF<sub>3</sub> THEN R7=H;  
 IF R6 = ALKOXY <1-2> AND R7 = H THEN  
     RESTRICT R6 = [3] ;  
 IF<2>R 6-7 = HALO / ALKOXY <1-2>  
 THEN BEGIN  
     RESTRICT R6 = [3] ;  
     RESTRICT R7 = [4,5]  
 END.

Figure 3.11

Had the expression been written without parentheses, it would have indicated an unsubstituted phenyl group, or a cyclohexyl group substituted by a methyl group, itself further substituted by chlorine.

This "highest-level" convention defines the level of substitution on substituents used as substituent values. In the statements

```
R1 = R2 sb Cl ;  
R2 = phenyl sb methyl
```

it is understood that the chlorine is a substituent on the phenyl group.

### 3.8. SPECIAL RESTRICTIONS IN GENERIC STRUCTURES

GENSAL provides two types of statement which allow special restrictions to be placed on the variety of specific structures covered by a generic structure : "IF" statements and "RESTRICT" statements.

Both employ the syntactic construct, condition, the former using the result of the condition (TRUE or FALSE) to determine which of two alternative statements should be used, and the latter to impose limitations on the definitions already made.

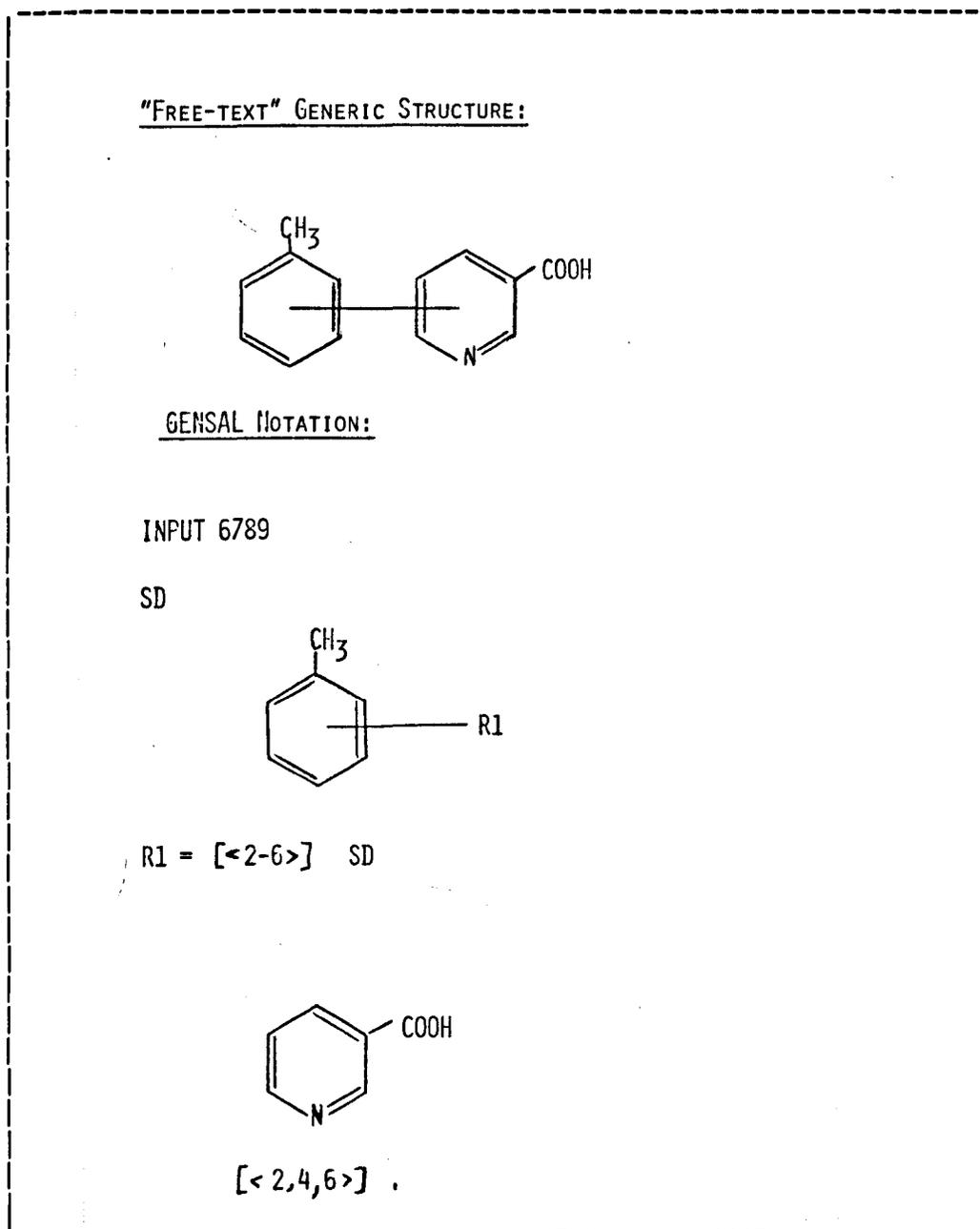


Figure 3.12

3.8.1. Conditions

Complex conditions can be formed, using the Boolean operators

AND, OR and NOT, as shown in Syntax Diagram 19. In executing such conditions the unary operator NOT is evaluated first, followed by AND, and finally OR; conditions in parentheses are evaluated first of all.

Ultimately, all conditions are composed of simple conditions of the form shown in Syntax Diagram 18. All simple conditions have two sides, separated by a relational operator ("=" or "<>", meaning "is" or "is not" respectively), and there are basically three types, which will be discussed separately. Each of them describes a particular arrangement of possible values for the variables in a generic structure, and for this reason only substituents and multipliers that have already been defined may appear in conditions.

### 3.8.2. Definition Relations

In "definition relations" the right-hand side is a substituent definition, of exactly the same form as is used in assignment statements, though here it may be abbreviated to a "stand-alone" position set, where the chemical nature of the substituent is not relevant. The left-hand side consists simply of a substituent or substituent combination, as in:

a) IF R1 = [4] methyl THEN...

(If R1 is a methyl group in the 4 position then...)

b) IF R2 + R3 = [2/3] THEN...

(If the attachments of the structure formed by the combination of R2 and R3 are at positions 2 and 3 then...)

c) IF R4 = alkyl<1-6> SB (Cl / Br / I) THEN...

(If R4 is an alkyl group of between one and six carbon atoms substituted by Cl, Br or I then...)

### 3.8.3. Integer Relations

Integer relations have a selector, identifying an integer range on the right-hand side, and the left-hand side can consist of various integer terms such as multipliers and substituents with parameters combined by arithmetic operators.

d) IF R1 C = <1-2> THEN...

(If the carbon count of the **homologous series identifier** defining R1 is in the range 1 to 2 then...)

e) IF R2 E = <2-> THEN...

(If there are two or more double bonds in the **homologous series identifier** defining R2 then...)

f) IF M1 = <2-3> THEN...

(If M1 is either 2 or 3 then...)

g) IF M1 + M2 + R1 C = <4> THEN...

(If the sum of M1 and M2 and the carbon count of R1 is 4 then...)

h) IF R1 C + R2 C = <12-> THEN...

(If the sum of the carbon counts of R1 and R2 is greater than or equal to 12 then...)

i) IF R1 + R2 C = <0-6> THEN...

(If the carbon count of the combined substituent formed by R1 and R2 is less than or equal to 6 then...)

The syntactic and semantic differences between the "+" symbols in examples (h) and (i) above are important. In the former it is an arithmetic operator combining separate integer values, whereas in the latter it combines the two substituents in a substituent combination.

3.8.4. Group Relations

"Group relations" begin with a selector which operates on the remainder of the left-hand side of the condition. If this is a substituent group, then the right-hand side will be a substituent definition, or stand-alone position set, as in the definition relations described above:

j) IF <2-> R1-5 <> H THEN...

(If two or more of the substituents R1, R2, R3, R4 and R5 are not hydrogen, then...)

On the other hand, if the left-hand side is a multiplier group or substituent group and parameter, then the right-hand side will be an integer range:

k) IF <1-3> M1-5 = <4> THEN...

(If 1, 2 or 3 or the multipliers M1, M2, M3, M4, and M5 is equal to 4 then...)

l) IF <1> R1 + R2 , R3 + R4 C = <3> THEN...

(If the carbon count of the group formed by either R1 and R2 or by R3 and R4 (i.e. if the carbon count of 1 of the two substituent combinations) is 3, then...)

3.8.5. IF Statements

In an IF statement, there are two subordinate statements after the condition, one following the delimiter THEN, and the other the delimiter ELSE (though this latter may be omitted). The statement following the THEN is used in those arrangements of the variables that make the condition TRUE, and the following the ELSE (if present) in those that make the condition FALSE.

The statements in the THEN and ELSE parts may be assignment statements, RESTRICT statements (described below), nested IF statements, "empty statements", or groups of statements enclosed within BEGIN and END delimiters (a "compound statement").

Examples of IF statements are:

```
IF R1 = methyl THEN R4 = methyl;
```

```
IF R1 = H THEN R2 = H
      ELSE R2 = halogen;
```

```
IF R1 = halogen
      THEN IF R1 = [2-3]
            THEN RESTRICT R2 = H;
            ELSE
      ELSE BEGIN
            RESTRICT R2 <> H;
            RESTRICT M1 = <3>
      END
```

There is no semicolon between the statement following the THEN and the ELSE, though the individual statements in a compound

statement are separated by semicolons.

In nested IF statements, each ELSE is paired with the most recent unpaired THEN: this can cause problems if there are nested IF statements without an ELSE part. In the last example above, an empty statement is used to provide an ELSE to pair with the second THEN, so that the effect is as intended. This point is further discussed in Section 3.11.1 below.

Clearly, certain IF statements would make semantic nonsense.  
e.g.:

```
IF R1 = methyl THEN R1 = ethyl
```

Such statements are illegal: if the condition involves a given substituent or multiplier then the statements in the THEN and ELSE parts may not involve that substituent or multiplier. The exception to this rule is that if the condition is concerned with the chemical nature of a substituent, then the statements may be concerned with its position, and vice versa. Thus the following statements are legal:

```
IF R1 = methyl THEN RESTRICT R1 = [2];
```

```
IF R2 = [4] THEN R2 = halogen
```

Figures 3.5, 3.8, 3.10 and 3.11 show examples of the use of IF statements from actual patent examples.

### 3.8.6. RESTRICT Statements

RESTRICT statements are used directly to reduce the possible arrangements of values for the variables in a generic structure. Only those arrangements which allow the condition to be TRUE are possible.

The form of the condition is exactly as in the IF statement, and thus RESTRICT statements appear as in the following examples:

a) RESTRICT R1 = H

H must have been given as a possible value for R1 in its original definition, and this statement eliminates all the other possibilities.

b) RESTRICT M1 + M2 <> <6>

It does not matter what the original definitions of M1 and M2 were; the RESTRICT removes those combinations of possibilities where their sum is 6.

c) RESTRICT R1 C = <2-3>

The carbon count of the **homologous series term** defining R1 is limited to 2 or 3 (which must be a subset of the values given in the original definition).

Examples from actual patents in which RESTRICT statements are used are shown in Figures 3.3, 3.8, 3.10 and 3.11.

### 3.9. SCOPE OF DEFINITIONS

There is no restriction on the number of different assignment statements each substituent or multiplier appears in, and all the different definitions given are alternative to each other.

When an assignment statement appears in the THEN or ELSE part of an IF statement, the alternatives given there for a substituent or multiplier are added to those given elsewhere (if any) when the condition has the appropriate value. If it is desired to limit the alternatives already given, then a RESTRICT statement should be used.

### 3.10. LIMITATIONS OF GENSAL

GENSAL has been designed to conform as closely as possible to the forms of expression commonly encountered in patent specifications and abstracts, whilst retaining a sufficient formalism for automatic processing to be possible. However it is not completely comprehensive, and at least in its present form, it is not applicable to certain types of expression found in patents.

In the majority of cases, the expression in question can be

reformulated in such a way as to permit encoding in GENSAL: the replacement of the different symbols used for structural and multiplicative variables by standard GENSAL substituents and multipliers is a trivial example of this. Other such limitations of GENSAL are the restrictions that structural variables may be at most doubly connected (Figure 3.11 shows how a small amount of respecifying of structural variables can circumvent this), and that only structural variables may have multipliers applied to them (again in Figure 3.11, this restriction is avoided by making what in the abstract is a multiplied structural constant, a structural variable with a selector applied to its single alternative value.

Certain other expressions found in patent specifications and abstracts cannot be represented in GENSAL at all, however, and several of the examples in the Figures show this.

The Derwent Abstract for the generic structure in Figure 3.7 indicates that certain of the alternatives are "preferred". At present, the only way around this problem is to construct a GENSAL notation for the structure in which only the preferred alternatives are shown; this could be stored alongside the more general notation. A fairly simple extension to GENSAL might allow the sequence of alternatives in a substituent definition expression to be interrupted by the delimiter PREFERABLY, those alternatives following it being the preferred ones. However, this could cause complications in a search system for generic structures encoded in GENSAL.

In the structure shown in Figure 3.10 it is not possible to show the limitations on R5 adequately using GENSAL. The requirement that it be "sterically hindered" cannot be shown at all (unless it be by indicating some branch points in the parameters), and that the majority of the occurrences of R5 must be alkyl<3-> can only be shown by exhaustively enumerating all the possible combinations of R2 and M1 in separate IF statements. This is reasonably satisfactory here, but would not be were there a much larger number of possibilities.

Figure 3.11 illustrates the lack of facilities to show stereochemistry, which is largely a consequence of the absence of stereochemical indicators in the two-dimensional structure representation used by the Feldmann graphics system. If GENSAL were to be used with a graphics system incorporating stereochemistry the syntax of GENSAL could be modified to include stereochemical descriptors in definition elements, treating them in a similar way to position sets.

Whilst GENSAL is not comprehensive, experience in encoding generic structures from patents suggests that in its present form it is capable of representing adequately the vast majority. However, the possibility is discussed in Chapter 6 that some modifications and extensions may need to be made to it.

### 3.11. THE DESIGN OF GENSAL

#### 3.11.1. Formal Grammar

The initial attempts at the design of GENSAL were based on analogy with Pascal, rather than on a rigorous approach using the formal grammar theory described in Chapter 2. No attempt has been made at formal proof of particular properties of the Grammar, but a number of such properties can be identified by intuitive inspection of the syntax diagrams shown in Appendix 1, or the equivalent Backus-Naur Form production rules shown in Appendix 2.

The Grammar of GENSAL is context-free, the production rules conforming to the requirements of Chomsky Type 2 Grammars (Section 2.1.1). It is unambiguous, and is a member of the class of LL(k) Grammars defined by Lewis and Stearns<sup>83</sup> (Section 2.2.2), which means that it can be parsed "top-down" as well as "bottom-up" (Section 2.2.3).

The syntax for IF statements in GENSAL is similar to that of Pascal and Algol 60, and as was pointed out in Section 3.8.5 can lead to difficulties with nested IF statements where not all have an ELSE part. This is one aspect of the design of Pascal which has been criticised.<sup>128</sup> Algol 68 and certain other languages have explicit terminators for IF statements ("FI" in Algol 68) which help to avoid this problem; this is an aspect of GENSAL

syntax which could perhaps usefully be modified.

### 3.11.2. Non-Determinacy

In most cases, inspection of the next symbol of a sentence in GENSAL is adequate to decide which production is being used: were this always the case the Grammar would be LL(1), and would be deterministic according to the definition of Koranjak and Hopcroft.<sup>86</sup> [The Grammar of Pascal is of this type.] However there are three places in the Grammar of GENSAL where lookahead is required for parsing.

In integer ranges (Syntax Diagram 2) there are three possible productions starting with integer:

<integer>

<integer> - <integer>

and <integer> -

and which of these is being used cannot be decided until up to two further symbols have been examined.

In position sets (Syntax Diagram 6) it is not possible to decide whether or not a position combination is being read until a plus sign is or is not encountered after the first integer.

In substituent groups (Syntax Diagram 12), both productions start with an R delimiter, followed by an integer, and only when the symbol after the integer is examined is it possible to decide

which is being used.

In simple conditions (Syntax Diagram 18) the situation is more complex. The possibility of a stand-alone position set instead of a full substituent definition expression on the right hand side of group definition relations and definition relations (Sections 3.8.2 and 3.8.4) means that which of the two is being used may not be decidable until the symbol following a position set is inspected, and as the position set may be of arbitrary length, no limit can be set on the amount of lookahead required.

Wirth <sup>109</sup> has suggested that where there are only a few examples of non-determinacy in the Grammar of a language, these should be handled on an ad hoc basis in the writing of a parser, and this approach has been adopted in the programming of the GENSAL interpreter, described in Chapter 5, where no particular difficulties were encountered with integer ranges and substituent groups.

The programming of the analysis of conditions has not formed part of the present work, and so the problem of arbitrary lookahead has not been considered in detail. However, there appears to be no need to know which path is being followed when a position set is encountered at the start of the right-hand side of a simple condition, and lookahead is therefore unnecessary. There is no reason to suppose that the semantic analysis of simple conditions would require this knowledge at the outset, and so it has been considered that the syntax for simple conditions is satisfactory

in its present form. In fact, though strictly speaking lookahead is required in the analysis of integer ranges, the interpreter program that has been written for GENSAL does not actually look ahead at all in performing their analysis, since it records the value of the first integer encountered, and later decides what to do with it.

### 3.11.3. Security vs. Flexibility

The relationship between the security and the flexibility of formal languages was discussed with reference to the programming languages Ada and Pascal in Sections 2.3.3 and 2.3.4. It was pointed out that the more redundant information is included in sentences in a language, the greater are the possibilities for the checking of self-consistency etc. The requirement in Pascal that all variables be declared with an indication of their type before they are used is an example of this; the type could in many cases be perfectly well deduced from the type of expression used in assignments to the variable in question.

GENSAL has however been designed to conform as closely as possible to the types of expression commonly found in patent specifications, and flexibility rather than security has been the principal aim. This is not to say that interpreter and compiler programs for GENSAL are likely to allow large numbers of errors to pass through undetected, but it makes the task of detecting and reporting such errors much more difficult.

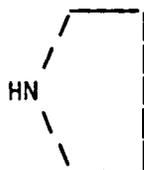
A more secure language for the description of generic structures might require all structural variables to be listed at the outset, with information on the number of connections and the bond orders for each. This information could then be used to check every occurrence of each variable. However this enhanced security would be at the cost of the natural form of expression currently found in GENSAL.

In GENSAL as it stands, information on the connectivity and bond orders of substituents frequently does not become available until well after the variable in question has been introduced, and in some circumstances may not become available at all, leaving the interpreter program to make assumptions about it. For example the GENSAL definition expression

phenyl sb R1

says nothing about the way R1 is connected to the phenyl group, or about its bond orders. When R1 is defined, it is possible that the information is still not given:

R1 = SD



In this case the program should assume that the connection is single, and that the bond order is single also, but a subsequent position set following the **structure diagram** in the manner described in Section 3.6 might give position combinations indicating that the connection was actually double.

Chapter 5 describes the approaches that have been used to detect incompatibilities in the information provided in a GENSAL sentence in the writing of an interpreter for GENSAL. In many cases the incompatibility can only be detected some time after an error has occurred, making recovery from the error very difficult if not impossible.

It is believed that the emphasis on flexibility rather than on security in GENSAL is justified, since most errors can be detected eventually, and it is GENSAL's flexibility, readability and similarity to the language of patent specifications and abstracts, in relation to alternative coding methods for generic structures, that is likely to be a major factor in determining its acceptability to the chemical and patent documentation industries.

## CHAPTER 4

### THE INTERNAL REPRESENTATION

"Look beneath the surface; let not the several  
quality of a thing nor its worth escape thee"

Marcus Aurelius Antoninus (121-180)

The last Chapter gave a description of the formal language GENSAL, which has been designed to encode generic chemical structures from patents in a form which can be processed by computer, yet which remains readily intelligible to a chemist or patent agent. The formalism of its Grammar makes it comparable to a high-level programming language, and thus the program which analyzes it can be thought of as equivalent to a compiler. To extend this analogy further, the internal representation of a generic structure which this program produces can be thought of as being equivalent to the object code produced by a programming language compiler, though unlike the object code for a

programming language, the internal representation is a machine-level data structure, rather than a set of machine-level instructions. Furthermore, as the analysis program is expected to operate interactively, it is better described as an interpreter than as a compiler.

In a generic structure information system, this internal representation can be used to generate fragments for use in searching, or directly for atom-by-atom tracing in the final stage of a search. In order to enable it to perform these functions satisfactorily, and yet remain in a form which can easily be generated from GENSAL input, a number of features have been incorporated into its design, and these will be described in this chapter.

#### 4.1. REQUIREMENTS FOR THE REPRESENTATION

Chapter 1 discussed the need for a full and unambiguous description of the generic structure, from which fragment screen descriptors of various types could be generated algorithmically, and the reasons for the selection of connection tables as the appropriate basis for this representation.

The purpose of the representation described here is not to store explicitly all the possible specific structures covered by a given generic structure, but rather to contain sufficient

information for exhaustive generation of all the specific structures to be possible, even though in most cases such an operation would be pointless, as well as computationally unfeasible where the number of specific structures covered is large, or even infinite.

Since the representation is to be built up from a generic structure input to the computer in GENSAL, the conversion problems are greatly simplified if certain features of the representation mirror features of GENSAL. In particular, as the syntax for the definition of substituents in GENSAL is essentially recursive, the structure of the internal representation should be recursive also.

GENSAL employs Geivandov's concept <sup>17</sup> of a generic structure as a (possibly vestigial) constant part, to which are attached variable parts which can vary in their chemical nature, position of attachment and multiplicity of occurrence, and which may themselves be further substituted by other constant and variable parts, down to any level. At each level, certain of the values for the variable parts may be alternative or additional to each other in complex nested Boolean relationships. This suggests two principal components for the internal representation, one containing information about the chemical nature of the constant and variable parts and the other information about the way in which they are connected together in terms of positions and multiplicity, and the Boolean relationships between them.

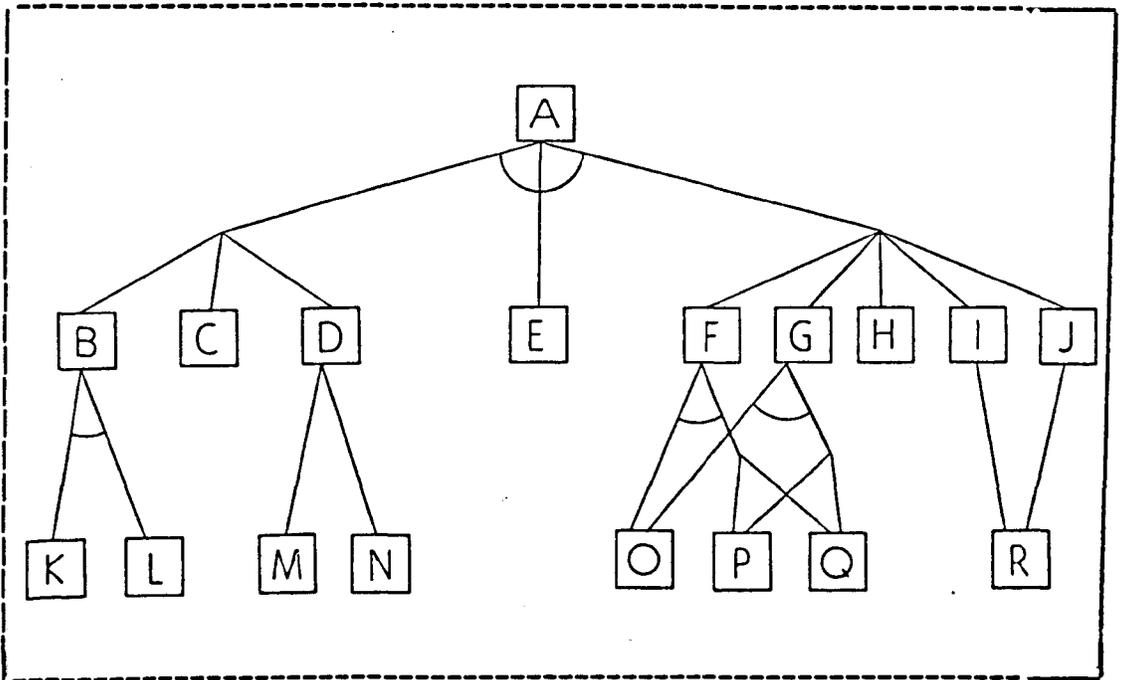


Figure 4.1: A diagrammatic representation of the basic structure of the ECTR showing the child gates. Each box represents a **partial structure**, and the lines represent **child gates**. Each hierarchical level of substitution is shown as a separate row of **partial structures**.

Lines meeting at a point connect together **partial structures** which are alternative to each other (OR relationship), and lines meeting at a point that are linked together by an arc connect **partial structures** which are additional to each other (AND relationship).

The GENSAL statements corresponding to this ECTR are shown in Figure 4.2.

The successive levels of further substitution imply a hierarchical relationship between the different parts of the structure, though the exact nature of the hierarchy depends on the way in which the GENSAL description of the structure was constructed, which is to a certain extent arbitrary. Where, as is illustrated in Section 3.7.4., a GENSAL substituent is defined in terms of itself, the hierarchy "loops back" to a higher level and there is no lowest level of substitution; the structure in

question is a polymer.

This approach to the storage of polymers has certain conceptual similarities to that developed by the Du Pont company in the 1960's.<sup>182</sup> In that system, each monomer unit is shown in a connection table connected to a dummy central atom, and path tracing procedures are able to pass through this central atom and back into the monomer unit(s).

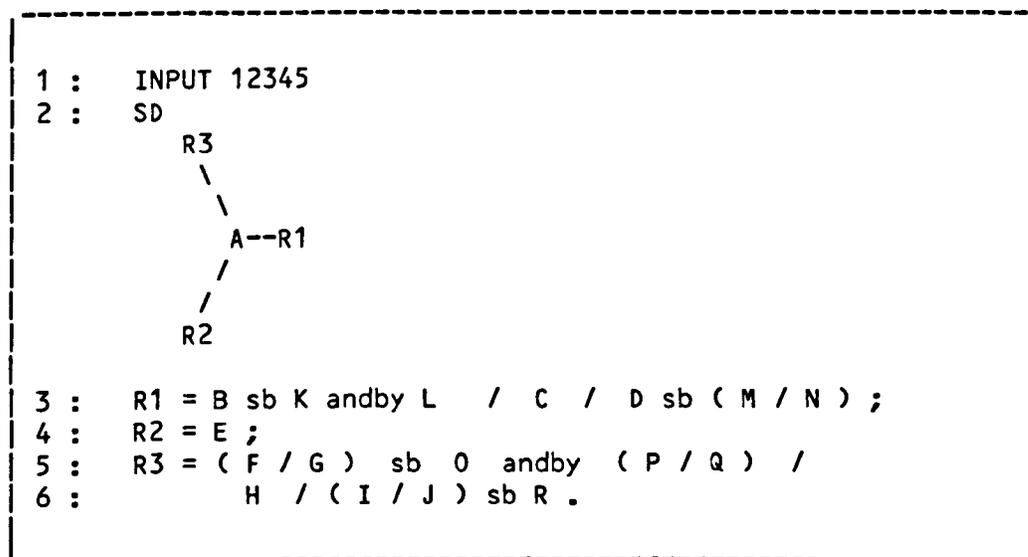


Figure 4.2: GENSAI statements corresponding to the ECTR's shown in Figures 4.1 and 4.4

Together the two components of the internal representation can be considered as forming a topological graph - the chemical nature of the various parts of the generic structure being represented in the nodes of the graph, and the information about their connections and relationships in its edges. Since information on the chemical nature of each part is predominantly based on conventional connection tables, the whole is a sort of super-

connection table, or connection table of connection tables, and is called an **Extended Connection Table Representation (ECTR)**. Within the ECTR each node is called a **partial structure (PS)**, and each edge a **gate**. The gates are divided into **child gates** and **parent gates**, according to which direction in the hierarchy they point; the graph is thus a directed one. The overall layout of the ECTR for a generic structure, showing the PSs and the **child gates**, is shown in Figure 4.1.

The entire ECTR is held in the main computer memory during its generation because, as further parts of the structure are defined during the course of the GENSAL sentence, it is frequently necessary to refer back to previously-defined parts. Similarly, as fragments are generated, or an atom-by-atom search performed, it is necessary to trace from one PS to another.

#### 4.2. THE PARTIAL STRUCTURE RECORD

Syntax Diagram 10 shows five different paths for a substituent value in GENSAL, and these were discussed in Section 3.3.1. From them it is possible to identify four fundamentally different types of **partial structure**, each requiring a different form of representation in the ECTR.

#### 4.2.1. Specific Partial Structures

These correspond to a single fully-defined structural entity, and are the only type of PS that may be represented by a connection table. They appear in GENSAL substituent values as **structure diagrams** (Section 3.3.1.2), or as **specific nomenclatural terms** (Section 3.3.1.3) which the GENSAL interpreter program translates into connection tables via a dictionary of standard nomenclatural terms.

#### 4.2.2. Generic Partial Structures

These appear in GENSAL substituent values as **homologous series terms** (Section 3.5), with associated **parameter lists**. They are shown in PS records as expanded **parameter lists**, including those parameters implied by the **homologous series term** itself, as well as those given explicitly. For example the term "alkenyl" implies at least one double bond, that would be indicated by the parameter E<1-> in the dictionary. This type of PS is designed to be handled for fragment generation and searching using the chemical grammars developed by Welford. 67, 174

4.2.3. Unknown Partial Structures

These appear in GENSAL substituent values as a "?" (Section 3.3.1.1.). Clearly, no further information can be stored about their chemical nature, and search algorithms should allow them to be matched against any structural entity.

4.2.4. Other Partial Structures

These cannot be directly associated with any particular structural characteristics, and include expressions such as "electron withdrawing group" or "easily hydrolysed group". They are shown in PS records as a character string, taken from the **other term** in the GENSAL substituent value, and could be used for some sort of text-based searching. Nomenclatural terms not found in the dictionary file used by the GENSAL interpreter program can also be stored in this form.

Table 4.1 summarises the information given in a **partial structure**.

Child Gate			
Parent Gate			
Specific	Generic	Unknown	Other
Connection Table	Parameter List	-	Character String

Table 4.1: Partial Structure Record

4.3. CONNECTION TABLE FORMAT

The connection table used to represent **Specific PSs** is a simple redundant one, each row representing one node, which may be either an atom (in which case the atom type is recorded as a two-letter symbol) or a GENSAL substituent (in which case its name - the "R1", "R2" etc of GENSAL - is recorded, along with the values it can take, in the same format as a **child gate**). The record structure for a connection table row is shown in Table 4.2.

Atom Type	Substituent Name
	Substituent Values (Child Gate)
Charge	
Number of Hydrogens	
Six Congeners	

Table 4.2: Connection Table Row

Normally, substituents attached to a **Specific PS** are not explicitly included in the connection table as information about the atoms to which they are connected is stored in the **child gates**. It is only when there is a chain (cyclic or acyclic) of such substituents connected together, as shown below, that it is necessary in order to indicate the order in which they are



set for each connection table, which is thus the maximum number of non-hydrogen atoms permitted in a structure diagram. However, because the splitting of a generic structure into separate PSs is to a certain extent arbitrary, a large structure diagram can always be divided into two or more smaller ones, and the limit might be different in other implementations.

None	Fraternal	Filial	Parental
-	Row number of connected atom or "NOTFIXED" for variable-position connection	-	-
Bond Order			

Table 4.3: Congener Record

#### 4.3.2. Bond Orders

In the present implementation, the different bond orders used have been derived from the Feldmann system,<sup>179</sup> with some modifications. Fifteen bond types are distinguished, and are shown in Table 4.4

Because the environment (chain or ring) of a particular bond may alter according to which alternative values for a particular structural variable are being considered, and because the possibilities for tautomerisation and aromaticity may change similarly, the finer distinctions between these bond types are

not always helpful in generic structures. Ideally, an operational system would permit the user only to distinguish between Single, Double, Triple and "Any" bonds, and would automatically perceive rings, tautomers and aromaticity. Algorithms for such analyses in specific structures have been developed for use in synthesis analysis programs. 183-185

Chain Single (CS)	Ring Single (RS)	Any Single (S)
Chain Double (CD)	Ring Double (RD)	Any Double (D)
Chain Triple (CT)	Ring Triple (RT)	Any Triple (T)
Chain Tautomeric (TC)	Ring Tautomeric (TR)	<i>Any Bond (A)</i>
Any Chain (C)	Any Ring (R)	
	Ring Alternating (RA)	

Table 4.4: Bond Orders in Connection Tables and Gates

#### 4.4. PARAMETER LIST FORMAT

Welford<sup>67, 174</sup> has described a means in which the parameters applied to a homologous series term in a GENSAL sentence can be used to apply constraints to the chemical grammars used for generation and/or recognition of the members of the homologous series. The standard parameter identifiers used in GENSAL to constrain such features as atom count, branch points and unsaturations are shown in Table 3.1, and substituents in

**parameter lists** can be used to indicate interruptions in a chain or ring, or substitutions on it.

The full set of parameters with their values is sufficient, when used to constrain the chemical grammars, to define completely all the possible structures covered. Consequently, the PS record for the **generic** type of PS can consist simply of a list of parameter values (as integer ranges) for all the standard parameters. The non-standard parameters, represented by GENSAL substituents, are treated as substitutions on the **generic PS**, and information about them is given in **child gates**, as described below. However, when generating fragments or path tracing within the ECTR, the information about **children** of **Generic PSs** can be used to apply constraints to the chemical grammars.

#### 4.5. CHILD GATE FORMAT

**Child gates** indicate the connections from one PS (called the **parent PS**) to those lower down in the hierarchy to which it is connected. There may be connections to several **child PSs**, which can be additional or alternative to each other. Each **child gate** therefore describes a "one-to-many" relationship, though over the ECTR as a whole the **child gates** between successive levels of the hierarchy describe a "many-to-many" relationship, as can be seen from Figure 4.1.

In order to show the Boolean relationships between the various

**child PSs**, as well as information on positions of attachment, bond orders etc., the internal structure of **child gates** is quite complicated. Each **child gate** is essentially a tree, with two different types of node; the root of the tree is attached to the **parent PS**, and the nodes are arranged in layers called **bars**. Each **bar** contains only one type of node, and is either a **combination bar** containing **combination bar item nodes**, which are in AND relationship, or is an **alternative bar** containing **alternative bar item nodes** in OR relationship. The two types of **bar** follow one another alternately.

For reasons of convenience, based on the precedence of operators in **GENSAL** expressions (Section 3.7), the information on positions, multiplicity, bond orders etc. is stored in the **combination bars**, which form the top and bottom **bars** of each **child gate**. The number of intervening layers depends upon the complexity of the Boolean relationships, as indicated by the number of pairs of parentheses in the **GENSAL** expression. It is possible for there to be only a single **combination bar** in a gate.

Both types of **bar** are constructed as linked lists of items, which are alternative to each other in **alternative bars**, and additional to each other in **combination bars**. The **child gate** field of a **PS** record (Table 4.1) is a pointer to the first item in the top **combination bar**, and each bottom **combination bar** points to a **child PS**.

Positions in Parent PS	
Multiplicity of Occurrence	
Bottombar	Not Bottombar
Positions in Child PS	Pointer to
Bond Order	Alternative Bar
Pointer to Child PS	(next layer)
Pointer to next item in Combination Bar list	

Table 4.5 : Item in combination bar of child gate.

4.5.1. Combination Bars

The record structure for a **combination bar** item, shown in Table 4.5, indicates that it may take one of two possible forms, according to whether or not it is located in the **bottom bar** of the **gate**. For both **bottombar** and **non-bottombar** forms, information is given about the positions of attachment in the **parent PS**, and the multiplicity of occurrence in these positions; there is also a pointer to the next item in the **combination bar**.

For **non-bottombar** items, there is a pointer to the first item in the **alternative bar** of the next layer, and the position and multiplicity information given applies to all the alternatives in this **alternative bar**.

For **bottombar** items, no such alternatives are possible, and a

pointer is given to the appropriate **child PS** record, along with information about the positions in the **child PS** at which the attachment may be made, and the order of the connecting bond.

Position information can be taken from explicit GENSAL position sets (if present) or calculated from those positions available for substitution; multiplicity information can be taken from a GENSAL selector or from the definition of a multiplier or, if the **child** has been specified in a **parameter list** for a **homologous series term** (Section 3.5), from the values given for that parameter.

If there are several **combination bars** in a **gate** then the position information may in each layer more closely specify the positions of attachment; the positions specified lowest down the **gate** are those that actually define the point of attachment in the **parent PS**. Not every layer necessarily has a value for the positions of attachment in the **parent PS**, but the top **bar** will always specify positions; others will only do so if there is a position set given in the GENSAL expression.

On the other hand, multiplicity information is given in every layer (and is assumed to be 1 if there is no other information), and the values in successive layers are multiplied together in the manner of Section 3.7.1.

4.5.2. Alternative Bars

These have a much simpler structure than **combination bars**, and the record structure for an **alternative bar** item is shown in Table 4.6. All the information about each alternative in the list is given in the **combination bar** pointed to.

Pointer to Combination Bar (next layer)
Pointer to next item in Alternative Bar list

Table 4.6: Item in **alternative bar** of **child gate**.

Figure 4.3 illustrates the internal structure of a single **child gate** for a moderately complicated GENSAL expression.

4.6. PARENT GATE FORMAT

The structure of **parent gates** is very much simpler than that of **child gates**, as none of the information on the Boolean relationships between the various **child PSs** is stored in them. In fact, all the information contained in a **parent gate** is also contained in the corresponding **child gates**, and the purpose of **parent gates** is simply to allow path tracing within the ECTR to

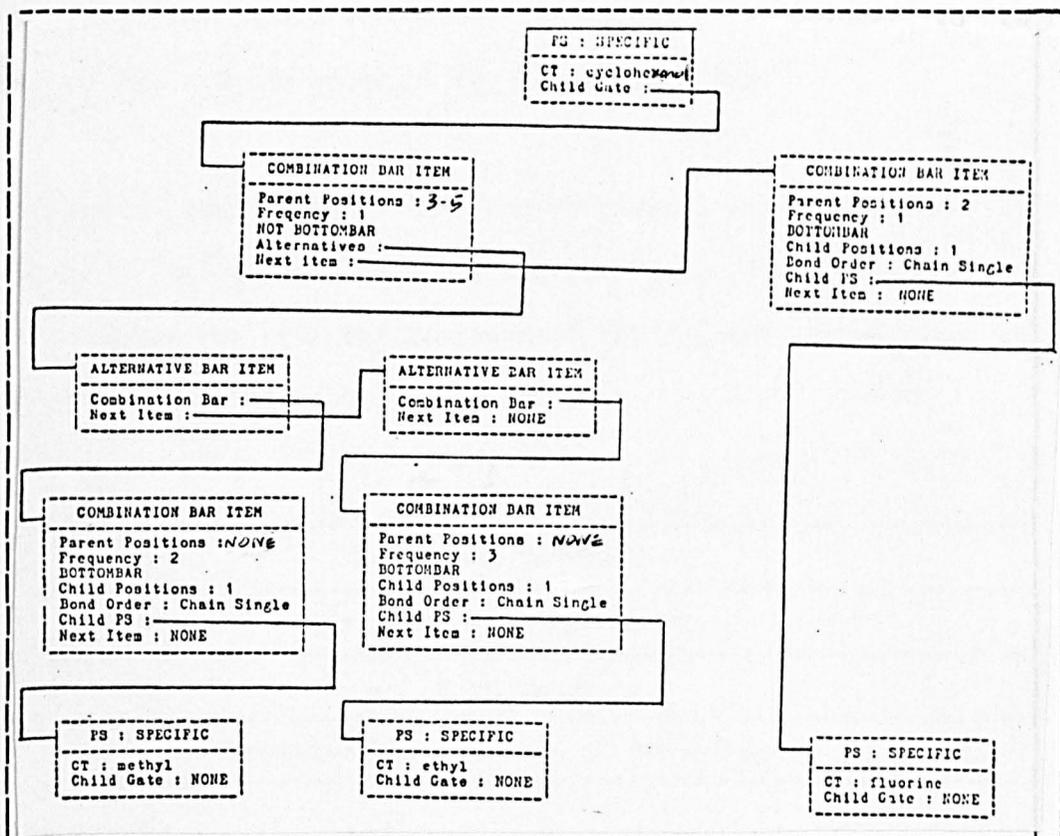


Figure 4.3: A diagrammatic representation of the structure of the child gate corresponding to the GENSAL expression: cyclohexanol SB [3-5] (<2> methyl / <3> ethyl) & [2] F which means that cyclohexanol is substituted in positions 3, 4, and/or 5 by either two methyl groups or three ethyl groups, and in addition to these by one fluorine in position 2.

take place from child PS to parent PS as well in the other direction; the redundancy of the information in the parent gates is compensated for by the substantial enhancements in path tracing ability.

Like the two types of bar in child gates, parent gates are implemented as a linked list of items, each item referring to a different possible parent PS for the child in question. The record structure is illustrated in Table 4.7. For each possible parent PS, the possible positions for connection in both the

child and the parent are given, along with a pointer to the parent PS, and the order of the connecting bond.

The parent gate field of a PS record gives a pointer to the first item in a linked list of parent gate items. Figure 4.4 illustrates the overall structure of the parent gates for the generic structure shown in Figure 4.2.

Positions in Child PS
Positions in Parent PS
Bond Order
Pointer to next item in Parent Gate

Table 4.7: Item in parent gate.

#### 4.7. REPRESENTATION OF CONDITIONS AND RESTRICTIONS

The ECTR described in this chapter makes no provision for incorporating the information given in GENSAL "IF" and "RESTRICT" statements, nor for distinguishing between the five different assignment operators that can be used to indicate independent or non-independent values for substituents or multipliers in selected group assignment statements (Section 3.4.2.).

These features of GENSAL, which mirror many of the expressions found in chemical patent specifications, are used to limit the

variety of possible specific compounds covered by a generic structure, by restricting the co-occurrence of particular alternatives in substituent definitions, etc. The present form of the ECTR may thus describe a greater variety of specific compounds than is actually warranted, and the limitations imposed by "IF" and "RESTRICT" statements could be implemented by indicating which of the possibilities in the ECTR should not co-occur. This might be achieved by applying some sort of selective "lock" to the gates, though the way in which this might be represented in the computer has yet to be determined.

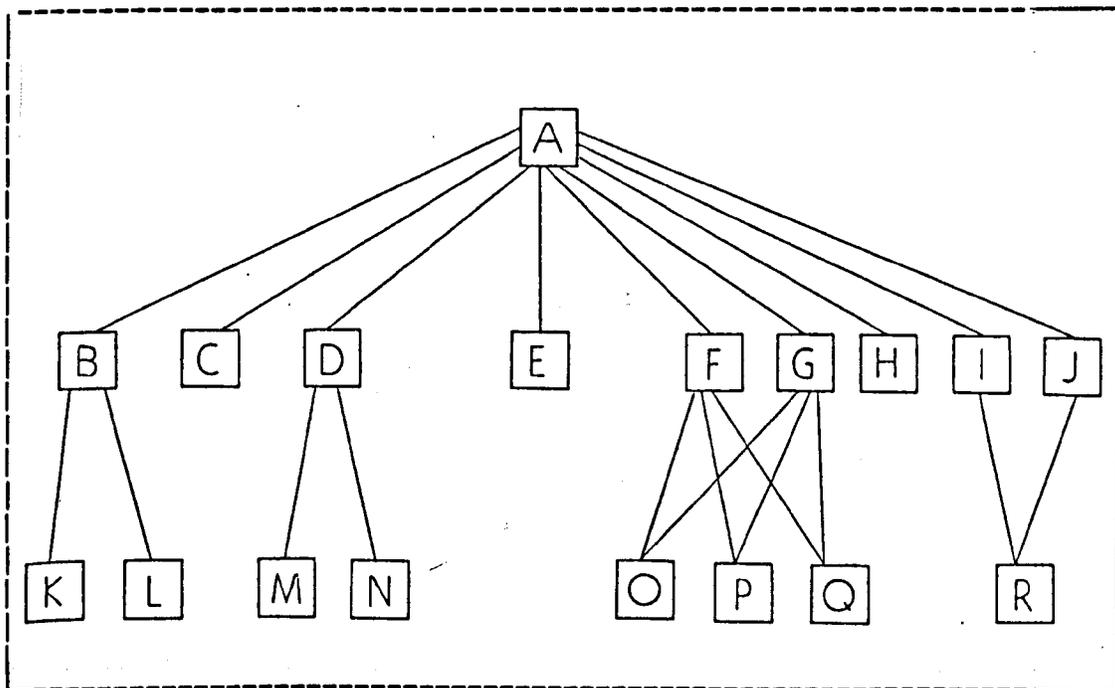


Figure 4.4: A diagrammatic representation of the ECTR, showing the parent gates, for the generic structure of Figure 4.2.

4.8. THE ECTR AND OTHER REPRESENTATIONS

Silk<sup>29</sup> has drawn attention to the similarity between a Markush structure and a nested Boolean expression, and suggested that the Boolean relationships could be incorporated into a notation-based representation for Generic structures. The ECTR, also exploits this similarity with the successive layers of bars in child gates representing the nested Boolean relationships, though the PSs are represented by connection tables, rather than notation strings.

An approach much closer to that described here has been proposed by Fugmann et. al.<sup>186</sup> It is based on an application to generic structures of the topological graphs used to represent concept relationships in the TOSAR (Topological Representation of Synthetic and Analytical Relations of Concepts) system developed by IDC.<sup>187</sup> Figure 4.5 shows the representation of a generic structure as a TOSAR graph which, like the structure of child gates in the ECTR, indicates AND and OR relationships between the different parts of the structure by means of two types of node in the graph (shown as open circles for OR and dots for AND). Fugmann et al. warn however, that the path tracing algorithms used for TOSAR graphs may be extremely expensive where tracing in generic structures is concerned.

The Chemical Abstracts Registry III System<sup>188</sup> employs a mechanism for compiling several partial connection tables to describe a larger specific structure. This involves the replacement of each ring system in a structure by a unique ring

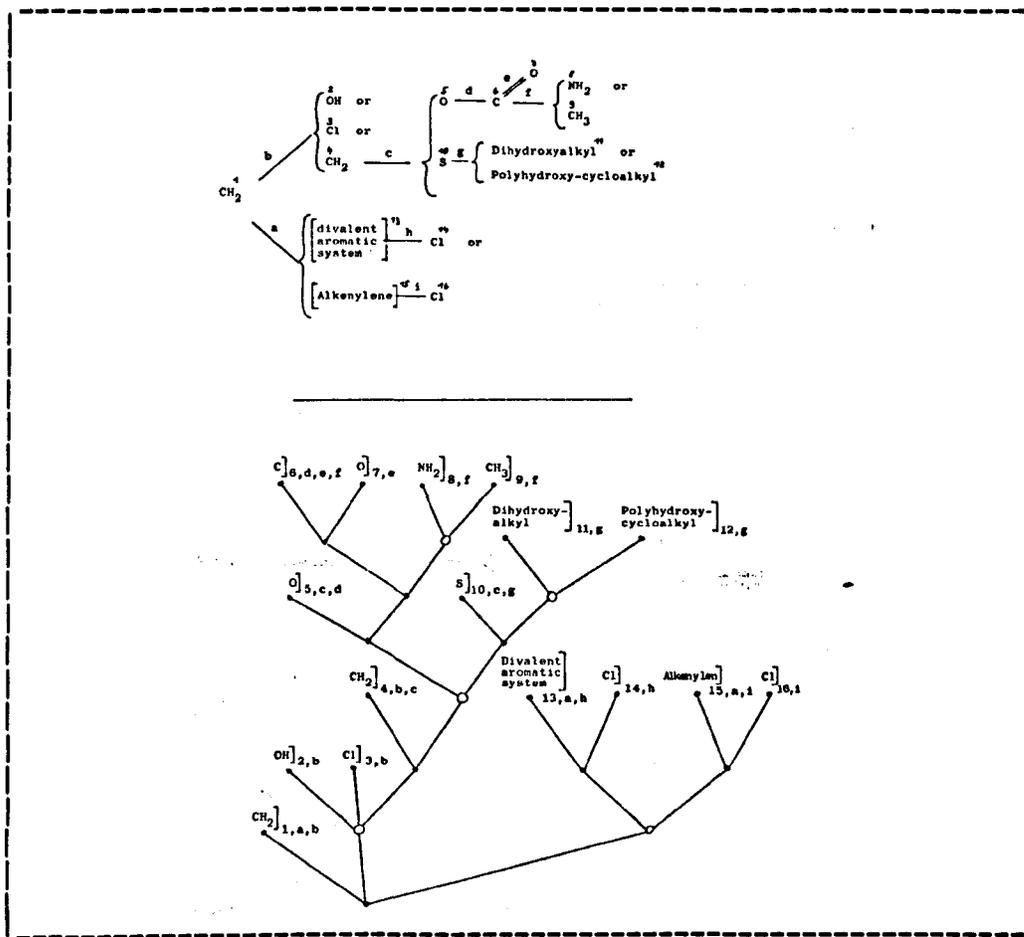


Figure 4.5: A generic structure represented as a TOSAR graph (From Fugmann et al.<sup>186</sup>).

identifier, which gives access to a separate file of connection tables for ring systems. This has the advantage of saving space, since only one connection table need be stored for each ring system, irrespective of the number of structures in which it occurs, and also allows cross-referencing between structures having ring systems in common, and acts as an aid to the automatic generation of systematic names<sup>189</sup> which are based on parent ring systems. The method is not used however as a means of describing generic structures.

#### 4.9. IMPLEMENTATION OF THE ECTR

The ECTR has been implemented using the data structures of the programming language Pascal. Because of the variable total size of the ECTR, which depends upon both the number and nature of the PSs, and its extensive use of linked lists, it is held entirely in dynamic storage, and access to its various parts is achieved using pointer variables.

Tenenbaum and Augenstein<sup>121</sup> have discussed the use of dynamic variables in Pascal, and more general problems of the implementation of recursive data structures have been considered by Hoare<sup>190</sup> and by Burton.<sup>191</sup>

##### 4.9.1. The Partial Structure Record

The Pascal TYPE declaration for a single PS record is a variant record, the variants corresponding to the four different types of PS found in the ECTR:

```

PTRPSTYPE = ^PSTYPE;
TPSVARIETY =(DUMMY, UNKNOWN, SPECIFIC, GENERIC, OTHER);
PSTYPE     = RECORD
    VISITED           : BOOLEAN;
    CHILDGATE        : PCOMBINLIST;
    PARENTGATE       : PPARENTLIST;
    CASE PSVARIETY   : TPSVARIETY OF
        DUMMY        : SUBSTNAME : SUBSTITUENT);
        UNKNOWN      : ();
        SPECIFIC     : (CT           : CTTYPER);
        GENERIC      : (PARAMLIST  : TPARAMLIST);
        OTHER        : (TERM       : STRING32)
    END;

```

In this record, the VISITED field can be used as an aid to path tracing in the ECTR, and the other two invariant fields give access to the child and parent gates respectively. Of the variant fields, that for a DUMMY PS is used only for housekeeping operations in the GENSAL interpreter program, no information can be stored for UNKNOWN PSs and the record TYPES for the other three varieties of PS are given below.

One of the advantages of using a variant record is that it is only necessary to set aside the amount of computer storage actually required for the particular type of PS in question.

#### 4.9.1.1. Connection Tables

The Pascal TYPE declarations are

```

CTTYPE     = ARRAY[1..MAXCT] OF ^ROW;
STRING2    = PACKED ARRAY[1..2] OF CHAR;

```

```

NUMCONGENERS=0..MAXCONGENERS;

SUBSTITUENT = 0..MAXVARS;

ROW          = RECORD
              CHARGE      : -9..9;
              HYDROGENS   : NUMCONGENERS;
              CONGENERS   : CONGARRAY;
              CASE ATOMICROW : BOOLEAN OF
                TRUE      : (ATOM      : STRING2);
                FALSE     : (NAME      : SUBSTITUENT;
                           VALUES   : PCOMBINLIST)
              END;

```

The connection table consists of an array of pointers to individual ROWs of the connection table; this is also a space-saving measure, as it means that there is no requirement to set aside large amounts of space to store empty connection table ROWs. MAXCT is a CONSTANT giving the maximum permissible number of ROWs, currently 32.

```

RELATIVES   =(NONE, FRATERNAL, PARENTAL, FILIAL);

ATOMNUMBER  = NOTFIXED..MAXCT;

CONGARRAY   = ARRAY [1..MAXCONGENERS] OF
              RECORD
              BOND       : BONDORDER;
              CASE RELATIONSHIP : RELATIVES OF
                NONE,
                PARENTAL,
                FILIAL    : ();
                FRATERNAL : (ROWNUM   : ATOMNUMBER)
              END;

```

In the array of congeners for each ROW in the connection table, the number of congeners permitted is controlled by the CONSTANT MAXCONGENERS, which is currently 6. A variant record distinguishes between the different types of connection, and in the ROWNUM recorded for FRATERNAL connections, a value of NOTFIXED (which is a CONSTANT equal to 0) indicates variable-position connection. The available

bond orders are

```
BONDORDER =(NOTSPECIFIED, ANY, CHAIN, RING, SINGLE, DOUBLE,
            TRIPLE, CHAISING, CHAIDOUB, CHAITRIP, CHAITAUT,
            RINGSING, RINGDOUB, RINGTRIP, AROMATIC, RINGTAUT);
```

#### 4.9.1.2. Parameter Lists

This consists an array of integer range records, one for each parameter, each consisting of a linked list of pairs of integers (being the lower and upper bounds of each subrange) plus a single integer to indicate the lower end of an unbounded top range:

```
PDOUBLIST  = ^DOUBLIST;

DOUBLIST   = RECORD
            FIRST,
            SECOND : INTEGER;
            NEXT   : PDOUBLIST
            END;

INTRECORD  = RECORD
            SUBRANGES : PDOUBLIST;
            TOPRANGE  : INTEGER
            END;
```

If there is no unbounded top range, then the TOPRANGE field is set to NOTSET, a CONSTANT of value -1

The declarations for the parameter list array are thus

```
TPARAMETERS =(ATOMCOUNT, TBRANCH, QBRANCH, EUNSATURATION,
              YUNSATURATION, RINGCOUNT, RINGATOMS,
              RINGSUBSTITUTION, RINGFUSIONS,
              RINGAROMATIC, HETEROATOM);

TPARAMLIST  = ARRAY[TPARAMETERS] OF INTRECORD;
```

## 4.9.1.3. Other Terms

This is simply a character string, currently of 32 characters:

```
STRING32    = PACKED ARRAY[1..32] OF CHAR;
```

4.9.2. Child Gate Record

The Pascal TYPE declarations for **combination** and **alternative bars** are:

```
PCOMBINLIST = ^COMBINLIST;

COMBINLIST  = RECORD
    PARENTPOSITIONS : PTGROUPMEMS;
    FREQUENCY       : INTRECORD;
    NEXT            : PCOMBINLIST;
    CASE BOTTOMBAR  : BOOLEAN OF
        TRUE : (CHILDPS      : PTRPSTYPE;
                CHILDPOSITIONS : TGROUPMEMS;
                CONNBONDS    : TCONNBONDS);
        FALSE: (ALTERNATIVES : PALTERNLIST)
    END;

PALTERNLIST = ^ALTERNLIST;

ALTERNLIST  = RECORD
    COMBINATION : PCOMBINLIST;
    NEXT        : PALTERNLIST
    END;
```

and they can be seen to correspond with the record formats shown in Tables 4.5 and 4.6. The *FREQUENCY* fields have the same TYPE as the elements of the parameter list record shown above, and the position set fields are as follows:

```
INTEGSET    = SET OF 0..MAXVARS;
```

```

TGROUPMEMS = RECORD
    CASE COMBINED      : BOOLEAN OF
        TRUE  : (COMBMEMS : PDOUBLIST);
        FALSE : (MEMBERS  : INTEGSET)
    END;

PTGROUPMEMS = ^TGROUPMEMS;

```

The **BOOLEAN** tag field for the variant record type **TGROUPMEMS** distinguishes between position sets for singly-connected substitution (**COMBINED = FALSE**), which are represented simply by an integer set, and position sets for doubly-connected substitution (**COMBINED = TRUE**), represented by a linked list of pairs of integers.

The **PARENTPOSITIONS** field of the **combination bar** item is a pointer to a **TGROUPMEMS** record, rather than a **TGROUPMEMS** record itself because, as was stated in Section 4.5.1., not all **combination bar** items have a record of positions in the **parent PS**, and for those that do not the **PARENTPOSITIONS** field can be set to **NIL**. Furthermore, the use of a pointer allows several different **combination bar** items to share the same **PARENTPOSITIONS^** record.

In contrast, there will always be information in the **CHILDPOSITIONS** field where **BOTTOMBAR** is **TRUE**, and thus this is a **TGROUPMEMS** record, and not a pointer to one.

The **CONN BONDS** field, showing the bond orders for the connection, is another variant record:

```

TCONNS      = NOTSET..2;

TCONNBONDS  = RECORD
  CASE CONNECTIONS : TCONNS OF
    NOTSET,
    0              : ();
    1              : (BOND : BONDORDER);
    2              : (BONDA,
                    BONDB : BONDORDER)
  END;

```

The tag-field indicates whether the substituent is unconnected (CONNECTIONS = 0), singly- or doubly-connected, an appropriate number of bond orders being given in each case. The NOTSET value for the tag field is used only in the setting up of the ECTR in the GENSAL interpreter, when it may not initially be known what the connections are.

#### 4.9.3. Parent Gate Record

This is implemented as a simple linked list of records, corresponding to Table 4.6:

```

PPARENTLIST = ^PARENTLIST;

PARENTLIST  = RECORD
  CHILDPOSITIONS,
  PARENTPOSITIONS : TGROUPMEMS;
  PARENTPS        : PTRPSTYPE;
  CONNBONDS      : TCONNBONDS;
  NEXT           : PPARENTLIST
  END;

```

#### 4.9.4. Space Requirements

As a complete and unambiguous representation of a generic structure, the ECTR is expensive in its storage requirements. For this reason, it is not intended that it should be stored permanently as a record of the structure. It would in any case be difficult to write the ECTR to a file and read it back into the computer on account of its complicated nature as a network of pointers.

It is expected that the ECTR would be built up during interactive input of a generic structure for a database of such structures, and then immediately used for the generation of fragment descriptors which would be stored for use in the first stages of searching. The ECTR would then be discarded, and could subsequently be regenerated from the stored GENSAL statements only if required for atom-by-atom matching in the final stage of a search.

The actual amount of core storage occupied by the ECTR depends, of course, on the size and complexity of the generic structure it represents. One containing a large number of different alternative values for a structural variable, all of which would have to be stored as separate PSs, would occupy much more space than one with only a few alternatives; the number of atoms in each connection table is also an important factor. The GENSAL interpreter program described in Chapter 5 is able to count up the amount of space being used, and the ECTRs for generic

structures from patents that have been processed by this program have ranged in size from 1156 to 10 674 PRIME 750 16-bit words. The Pascal implementation used for the interpreter<sup>138</sup> allows 16 segments of 64 kwords each for the storage of dynamic variables, making a total of over one million 16-bit words available, though other implementations might not be so generous.

## CHAPTER 5

### AN INTERPRETER FOR GENSAL

---

"This is the interpretation of the thing"

Daniel, Ch. 5, Vs. 26

This Chapter describes an interpreter program, written in the Pascal language, which implements a subset of the GENSAL generic structure description language, and which is upwards compatible with the full language, as described in Chapter 3.

The interpreter program performs syntactic and semantic analysis on sentences in GENSAL, and generates an Extended Connection Table Representation (ECTR) of the structure described.

It is implemented as a separately-compiled procedure of a program

called GENPROG, which is a prototype generic structure storage and retrieval system under joint development by the author and Welford. Appendix 3 is a listing of the interpreter, **procedure INTERPRET**, and Appendix 4 contains a line-number index to the subordinate **procedures** and **functions** within it. Appendix 5 is a listing of the **const**, **type** and **var** declarations that are global to GENPROG, with the addition of those **procedures** and **functions** called by INTERPRET which are also called by other parts of GENPROG.

Pascal programs are sufficiently clear to be largely self-documenting; comments at the start of each **procedure** and **function** indicate the routine's basic purpose, and list the calls to it. This Chapter gives an overall view of the strategies involved in the analysis of GENSAL sentences, and the build-up of the ECTR, with particular notes on the capabilities and limitations of the interpreter, and on the error messages given by the program. It is not intended by itself to give a complete understanding of the workings of the program, for which it should be read in conjunction with a thorough study of the program listings in Appendices 3 and 5.

Appendix 6 shows a sample interpreter session, illustrating the input of a generic structure from a patent.

### 5.1. INVOCATION OF THE INTERPRETER

The main part of GENPROG processes a simple command language which allows the user to invoke the interpreter, and also to perform a variety of other functions. These include filing and retrieving of structures processed by the interpreter, opening and closing of files of diagnostic information on the program, adding to a dictionary of nomenclatural terms and invoking a simple interactive editor program for structures encoded in GENSAL, which has been written by Kinsella.<sup>192</sup> Ultimately it is expected also to have facilities for searching a database of generic structures, using GENSAL-encoded query structures, and printing search results in a variety of formats.

The interpreter may be invoked in one of two modes: interactive mode, in which each new line of GENSAL is typed at the terminal, and non-interactive mode in which previously stored lines of GENSAL are processed. Such lines might have been stored in a file after a previous session, or be the result of editing a structure.

The lines of GENSAL are stored as a linked list of lines, with pointers to both the preceding and following lines; connection tables, representing structure diagrams within GENSAL, are encoded so that they may also be stored as character strings, as discussed in Section 5.5.3. below. The forward and backward pointers in the linked list are intended to facilitate operations in the editor module of GENPROG.<sup>192</sup>

## 5.2. LEXICAL ANALYSIS

This is the first stage of analysis in any compiler or interpreter <sup>80, 81</sup> and is the process by which the input string of characters is divided up into tokens, each representing one terminal symbol. In INTERPRET, the variable TOKEN holds the most recently-identified token for examination by the syntax analysis routines, and it is updated by the procedures NEXTTOKEN and LOOKAHEAD, both of which call procedure GETTOKEN, the lexical analyser itself.

Three different types of token are identified: GENSAL delimiter words and symbols, nomenclatural terms, and integers; the subordinate procedures and functions in GETTOKEN determine which of these is present. This is done by moving the pointer N along the global variable BUFFER, which contains an upper-cased version of the last line read.

This arrangement means that lower-case letters may be used in the input, but they are treated as if they were upper-case; the user may adopt his own conventions as to the use of lower-case letters for nomenclatural terms, or delimiter words etc. In addition, each line of input may be edited using backspacing etc. before it is processed. In the Pascal implementation used <sup>138</sup> a non-standard extension to the standard procedure READLN allows entire packed arrays of char to be read in a single operation, the right-hand end of the array being space-filled if necessary, and an extra variable returning the number of characters actually

read.

When the end of the line is reached, procedure READLINE obtains a new one from the terminal, adding it to the linked list of lines, if the interpreter is operating interactively, or obtains it from the existing linked list, if the interpreter is operating non-interactively. If there are no more lines in the linked list to be read, then the interpreter automatically swops to interactive mode, and in interactive mode, the user is able to exit from the interpreter by entering a blank line.

### 5.3. SYNTAX ANALYSIS

The basic approach used for syntax analysis is that of top-down, recursive-descent parsing, as described by Wirth.<sup>118</sup> No single part of procedure INTERPRET is entirely concerned with syntax analysis, since the procedures and functions which carry it out are also concerned with semantic analysis and ECTR generation.

The analysis of structure description (Syntax Diagram 21) takes place in the body of procedure INTERPRET, and the analysis of statements (Syntax Diagram 20) in procedure STATEMENT. Separate procedures exist for the analysis of assignment statements (Syntax Diagram 17), RESTRICT statements, IF statements and compound statements, the last two being of necessity mutually-recursive with procedure STATEMENT.

A group of nested **procedures** analyse substituent definition expressions (Syntax Diagram 15), these being **procedure** ALTNVLIST, which encloses **procedures** ALTNTVE and ELEMENT (analysing a definition element (Syntax Diagram 14)) which recursively calls **procedure** ALTNVLIST.

Conditions (Syntax Diagram 19) have not been implemented as part of the present work, and **procedure** CONDITION simply accepts any sequence of tokens until an appropriate terminator is encountered. This means that, whilst IF and RESTRICT statements are not actually implemented, no errors are generated by their inclusion. A **boolean** flag, CONDITIONSPRESENT, controls the printing of a warning message at the end of structures containing conditions.

A number of other **procedures** carry out syntax analysis on particular syntactic constructs in GENSAL. These are **procedures** INTEGERRANGE, SELECTOR, POSITIONSET, PARAMETERLIST and SUBSTGROUP.

#### 5.4. ERROR HANDLING

The program detects four different types of error, printing appropriate messages at the user terminal. In each case the error message required is obtained from an external file, ERRORMSGS and is printed by **procedure** WRITEMESSAGE, with the possible inclusion of some information on bond orders, atom numbers etc., if

relevant. The available messages are listed in Appendix 7, and the sample session shown in Appendix 6 illustrates several of them.

#### 5.4.1. Program Errors

A limited number of checks are performed by the program on its own working, and any error detected causes the user to be ejected from the interpreter, with display of a unique error number.

#### 5.4.2. Structure Diagram Errors

These are errors detected during the processing of structure diagrams, and relate to such matters as illegal valencies etc. If any are detected, the structure diagram is rejected and the user required to correct it before processing can continue.

Structure diagram processing is more fully described in Section 5.5 below.

#### 5.4.3. "Immediate" Errors

These are errors relating to invalid tokens in the GENSAL input, and they are handled by procedure ERROR. All syntax errors fall into this class, as do certain semantic errors.

A line of arrows is drawn under the offending token in the GENSAL input line currently being processed, followed by the error message. The remainder of the input line is ignored, and in interactive mode the user is invited to continue the GENSAL input starting with a replacement for the erroneous token. In non-interactive mode the user is ejected from the interpreter.

#### 5.4.4. "Delayed" Errors ("Failures")

This type of error is not detected until processing has continued for some time after the token which causes it has been obtained by the lexical analyser, and it is called a failure. Failures relate to such matters as incompatible bond types, and are handled by procedure FAILURE. In all circumstances the user is ejected from the interpreter.

#### 5.5. STRUCTURE DIAGRAM PROCESSING

As was stated in Chapter 3, the graphics system used for the input of structure diagrams in GENSAL is intended to be implementation-dependent, and in the present work a modification of the structure generation and display program written by Feldmann and others<sup>179</sup> is being used. This consists of some 4000 lines of Fortran, and is implemented as an EXTERNAL procedure of GENPROG. In order to avoid the complexities of attempting to link the COMMON blocks used by the Feldmann program for storage

of the connection table it uses with the global variables of GENPROG, the connection table is transferred to and from the Feldmann program via a scratch disc file.

The connectivity and bonding tables in the Feldmann program are separate, and are read into the GENPROG global arrays FELDCT and FELDBD by procedure READFELDMANN. Procedure PROCESSCT then reformats them into the connection table format used by the ECTR and described in Section 4.9.1.1.

#### 5.5.1. The Feldmann Program

The principal modifications made to the Feldmann program have been to allow the identification of a node in the diagram as a GENSAL substituent, or as an "apical" connection (\*) or as a "variable-position" connection (#) as well as as an atom of a particular element, and to allow multipliers to be applied to a particular node.

In addition to this, the maximum number of nodes permitted has been reduced from 100 to 32, and upon exiting from the Feldmann program all "default" bonds are replaced by either chain or ring single bonds, depending upon their environment.

Some slight changes have also been made to the bond types permitted, and to the symbols used to represent them in the diagrams, and routines have been written to output the

connectivity and bonding tables to a scratch file, and to read them back again.

### 5.5.2. Procedure PROCESSCT

This procedure is applicable only to the Feldmann graphics system, but is virtually the only routine in the interpreter to be so, and thus is the only one that would require replacement were a different graphics system to be incorporated.

Since the Feldmann program carries out very few checks on atom valencies etc., such checks are done by PROCESSCT, which uses procedure REJECT to handle any errors detected.

The Feldmann connectivity table in FELDCT is examined line by line, but only nodes representing atoms (except hydrogen) and certain substituents are added to the ECTR-format connection table (Section 4.9.1.1). Procedure HNUMBER is able to calculate the number of hydrogens (equivalent to positions available for further substitution) on each atom for common elements, obtaining the permissible valencies from an external file, VALENCYFILE.

The bond orders are represented by an enumerated type which is so arranged that the ordinal values correspond to the integers used for the bond types in the Feldmann program. Since "default" bonds are removed from the structure diagram, NOTSPECIFIED bonds cannot appear in connection tables.

PROCESSCT checks that the connectivity of each substituent is compatible with its connectivity in any previous appearances, rejecting the **structure diagram** if it is not.

If a **structure diagram** is rejected, the lines of the ECTR connection table are DISPOSEd, and in interactive mode the user is returned to the Feldmann program to correct it; in non-interactive mode, the user is ejected from the interpreter. Otherwise, procedure GETPOSNS is used to determine the connectivity, bond order(s) and possible position(s) of attachment of each substituent in the diagram, removing from the connection table those substituents that are attached only to atoms.

### 5.5.3. Storage of GENSAL Structure Diagrams

In order to allow the **structure diagrams** occurring in GENSAL sentences to be held in the same format as text lines of GENSAL, the Feldmann-format connection table is encoded as a character string (using the Pascal CHR function for the integers in the connection table). The conversion is carried out by the **procedures** ENCODECT and DECODECT.

Since it is the Feldmann-format connection table which is used for this, the lines of GENSAL stored in files etc. include **structure diagrams** in Feldmann connection table format. A possible minor enhancement to the program would be to remove this

Feldmann dependency, and thus make it easier to use the interpreter program with other structure graphics systems. This could however leave problems with graphics systems having differing requirements for the storage of 2-dimensional atomic co-ordinate data; the Feldmann program retains no such information, but recalculates co-ordinates every time the diagram is redrawn.

### 5.6. SUBSTITUENT DECLARATIONS

The program maintains a record of the substituents declared (introduced) and defined during the course of a GENSAL sentence. This allows it to check firstly that all declared substituents are defined somewhere (procedure CHECKALLDONE), secondly that only declared substituents are defined, and thirdly that all declarations of a given substituent are compatible in matters of connectivity and bond order(s). GENSAL substituents can be declared in one of four ways:

- (a) in structure diagrams
- (b) as a user-defined parameter to a homologous series term
- (c) as a value in the definition of another substituent
- (d) in copying a definition containing a declaration as in (c) above. (This last is internal to the program, and not apparent to the user.)

In each case, an entry is made by procedure DECLARESUBST in a table of substituent declarations, RDECLARATIONTABLE, which is an

array of linked lists, one for each **substituent**. Each new declaration of a **substituent** is recorded as a new item in the appropriate list, which contains information about the declaration relating to such matters as the **partial structure** in which it occurs, the position(s) at which the substituent can be attached and the order(s) of the connecting bond(s).

If it is found that a substituent being declared has already been defined, then the values with which it was defined are copied into the **child gate** of the **partial structure** in which the new declaration occurs. This is done by the mutually-recursive procedures COPYCOMBAR and COPYALTBAR which copy bars of **child gates**. "Absolute" definitions of each substituent are held in the elements of an array called RDEFINITIONTABLE, in order that the definitions copied are independent of the environment (positions of attachment etc.) in which the substituent in question had previously appeared.

Not all of the information for entries in RDECLARATIONTABLE is available at the time the declaration is made, and missing items are filled in later.

Where one substituent is defined in terms of another, as in

R1 = R2 sb methyl

there may be further substitution to attach to the substituent given as a substituent value. Because, if this new substituent

has not yet been defined, no **partial structure** exists to which a **child gate** can be attached, a **DUMMY partial structure** is created to represent the substituent, and the **FURTHERSUB** field of the entry in **RDECLARATIONTABLE** points to this **DUMMY partial structure**. When the substituent in question is defined, the further substitutions on it can be copied onto the **partial structures** representing its possible values.

### 5.7. SUBSTITUENT DEFINITIONS AND ECTR GENERATION

When a substituent group has been read, procedure **POINTERLIST** sets up a linked list, each item of which represents one **RDECLARATIONTABLE** entry for one substituent or substituent combination in the substituent group (plus one extra item for **RDEFINITIONTABLE**). This list is passed as a parameter (**PARENTPSLIST**) to procedure **ALTNVLIST**, which creates **alternative bars** in **child gates**, one **child gate** being built up on each of the items in **PARENTPSLIST**.

Procedure **ALTNVLIST** contains an iterative **repeat** loop which cycles round all the alternatives (separated by "/" delimiters) in a **GENSAL** substituent definition expression, and calls procedure **ALTNTVE** to analyse the definition elements separated by "&", "OSB" and "SB" delimiters. The **PARENTPSLIST** linked list is slightly reformatted before being passed as a parameter (**PARALTLIST**) to **ALTNTVE**, which passes it on to procedure **ELEMENT**, which analyses a single definition element and builds up the bulk

of the ECTR.

### 5.7.1. Syntactic and Semantic Analysis in ELEMENT

**Procedure** GETLIMITPOSITIONS is used to determine the set of positions available in the **parent partial structures** for all the items in PARALTLIST, which also contains information on position sets given in previous recursions of definition element. Thus **procedure** POSITIONSET is able to give appropriate error messages if illegal positions are specified.

No such checking is performed in the analysis of selectors in definition elements. Thus no error would be detected in the following expression

R1 = phenyl sb [2] <5> methyl

There is no reason in principle why such checking should not be done, but it would involve considerable computational effort, and it has not been considered worthwhile as the the interpreter is not intended to be a teaching program. For similar reasons, no error is reported in the analysis of statements such as

R1 = phenyl sb [2] (F & Cl & Br & I)

### 5.7.2. Substituent Values

The analysis of substituent values is performed in a case statement, with separate procedures to handle each path.

For parenthesised substituent definition expressions, function NEWPARENTPSLIST sets up a new linked list, based on the items in PARALTLIST, for passing in a recursive call to procedure ALTNVLIST. This function also adds an extra non-BOTTOMBAR combination bar to the various child gates accessed via the items of PARALTLIST.

Since GENSAL treats substituents occurring as substituent values as parenthesised expressions (Section 3.7.3), an extra non-BOTTOMBAR combination bar is incorporated into the child gates with the DUMMY partial structure created to represent the substituent (Section 5.6) being included as one of the ALTERNATIVES leading from it. This is done by function EXTRALAYER.

Structure diagrams, always preceded by the delimiter "SD", are handled by calls to the Feldmann program and procedure PROCESSCT; appropriate partial structure records are also set up for "?" and "other term" substituent values.

### 5.7.3. Nomenclatural Terms

The analysis of nomenclatural terms is quite complicated, and is handled by procedure TRANSLATENOMEN. The approach used in this implementation of GENSAL has been to maintain a dictionary of nomenclatural terms (SPSDICT) which gives access to a file of structure records (SPSFILE). The entries in SPSFILE may be of three types: a connection table, a set of homologous series term parameters, or a GENSAL expression. In order to allow synonyms to be handled, several different records in SPSDICT may give access to the same record in SPSFILE.

Function RECORDHELD determines whether or not a record is held for a particular nomenclatural term; if none is, then the term is treated as an "other term" and an OTHER partial structure is used to store the character string itself.

If there is a record, function SPSVARIETY determines which of the three possible types it is. Both SPECIFIC (connection table) and GENERIC (parameter list) entries can be handled quite simply. OTHER (GENSAL expression) entries are more complicated.

This type of SPSFILE record is used for compound nomenclatural terms, which can be analysed into simpler terms: examples include "halophenyl", "diethylamino" and "N-methyl-2-propionamido". Other such terms represent a delimited series of alternatives, such as "halogen" or "alkali metal". The SPSFILE entry for halophenyl gives the expression "phenyl sb halogen", and that for halogen

the expression "F / Cl / Br / I". The entries for phenyl, F, Cl, Br and I are all partial connection tables.

The processing of such an expression involves saving the current input BUFFER etc., and then calling ALTNVLIST recursively to analyse it; after the return from ALTNVLIST the original input BUFFER is restored. Such nesting of expressions can continue to any level, and the interpreter effectively treats each expression obtained from SPSFILE as if it were in parentheses (in fact the expression as it appears in SPSFILE is always terminated by a parenthesis).

The "highest-level" convention for further substitution on parenthesised expressions (Section 3.7.4) is of particular importance when dealing with compound nomenclatural terms and was chosen in preference to the alternative "lowest-level" convention on account of the problems that the latter would cause with such expressions.

If the expression "halophenyl sb methyl" were to occur in a GENSAL sentence, the dictionary-lookup operation would result in its being treated as

(phenyl sb ( F / Cl / Br / I )) sb methyl

and the highest-level convention means that the methyl group is attached to the phenyl group and not to any of the halogens.

The process of dictionary lookup effectively changes the "right-rooted tree" of the compound nomenclatural term (where the

rightmost part of the term is connected back to the parent structure) to the "left-rooted tree" of a GENSAL expression (where the leftmost part is connected back to the parent structure).

A minor problem remains with compound terms such as "alkoxy" which, if interpreted as

(oxy sb alkyl)

would imply in the expression

alkoxy sb chlorine

that the chlorine was substituted on the oxy group. This is not the generally-understood meaning of such expressions, and the problem is really a result of the conflict between the use of right-rooted and left-rooted trees in standard chemical nomenclature and GENSAL expressions (which derives them from the forms of statement in patents) respectively.

The interpreter program gets round the problem by a "fiddle" of dubious chemical validity and the SPSFILE entry for "alkoxy" is

alkyl sb [0/1] oxy

Thus the oxy group is regarded as being a child of the alkyl group rather than vice-versa, and the position combination [0/1] is used to indicate that the oxy group is interposed between the alkyl group and its parent structure.

This approach is justified firstly because it avoids a tricky problem, and secondly because it will allow the SPECIFIC oxy group partial structure to be handled along with the GENERIC

**partial structure** for the alkyl group in fragment generation. <sup>67</sup>  
Had the alkyl group been a **child** of the oxy, this would have been more difficult. In any case, so far as the GENSAL user typing "alkyl sb chlorine" is concerned, the whole arrangement is hidden, and he need not be aware of the construction of the ECTR.

It is possible that a compound nomenclatural term may be converted via SPSFILE to an expression involving a simple **homologous series term**; the terms "chloroalkyl" and "alkoxy" are examples of this. Any parameter list given after the term will thus be used to specify parameters for the simple **homologous series term** in the GENSAL expression. However, as this term may be nested in several layers of GENSAL expressions, the variable INSERTHSTPS is used to keep a note of any **GENERIC partial structure** in the expressions obtained from SPSFILE. Should more than one **homologous series term** be encountered during the processing of an expression from SPSFILE, it would not be clear which of them should be qualified by the parameter list; for this reason a program error is given in this situation, which should not arise if care is taken in the construction of SPSFILE.

Routines exist in the main part of GENPROG for adding records to SPSDICT and SPSFILE, though one of the problems in building up these files has been deciding how to interpret certain terms. This is a matter discussed by Dyson <sup>54</sup> and referred to in Section 1.4.6. For example, it is not always clear if the term "alkenyl" indicates exactly one double bond, or a minimum of one. Clearly, a decision of some sort has to be made for the purposes of

SPSFILE, but it is possible that an operational system might allow the user to redefine certain terms for his/her own use, and perhaps to maintain a private dictionary file.

Ultimately, the real problem is that the meaning of a term may differ from patent to patent, and may be left deliberately vague; sometimes patents define the meaning of a particular term used, but the definitions of a term like "aryl" differ widely from patent to patent. There appears to be no simple solution to this difficulty, which will only finally be overcome if the drafters of patents agree on standard meanings for the terms they use.

#### 5.7.4. Parameter Lists

Procedure PARAMETERLIST carries out the analysis of parameter lists, and checks that the values given for each parameter are a subset of those implied in the **homologous series term** to which the list is being applied. Thus, for example, any value other than 0 for the number of rings in a parameter list applied to the term "alkyl" would be illegal.

Since it is possible for an **homologous series term** to be missing from SPSDICT, procedure TRANSLATENOMEN permits terms not found to be followed by parameter lists, though the information they give is not stored in the ECTR.

### 5.7.5. ECTR Generation

Two procedures handle the creation of child and parent gates, respectively SETCOMBARS (which calls function NEWCOMBAR) and SETPARENTGATE.

SETCOMBARS uses procedure GETCHILDPOSITIONS to determine the positions available in the child structure for the connection(s) to its parent. This procedure additionally checks that the bond orders specified are compatible. For each connection the bond order may have been specified in both the parent partial structure and the child partial structure (though in many cases either or both of these will be NOTSPECIFIED). Procedure BONDHECK uses a table of bond orders, BONDMATCHARRAY, to determine a bond order compatible with the two specified: for example a CHAIN bond and a SINGLE bond result in a CHAISING bond, whereas a CHAISING and a CHAIDOUB bond are recognised as incompatible. Two NOTSPECIFIED bonds result in a SINGLE bond, so that no NOTSPECIFIED bonds are left when the ECTR is complete.

The child positions determined may be modified if there is a position set following the substituent value, this modification being achieved by procedure MODIFYCHILDPOSITIONS.

Where further substitution has been specified on substituents given as substituent values, as in

R1 = R2 sb methyl

**Procedure** ADDFURTHERSUBTN is used to copy the **partial structures** for this further substitution (methyl in the above example) onto the **partial structures** created when the substituent in question (R2 in the above example) is defined. This uses **function** PPOSNS to check that any position sets specified are actually available, before calling COPYCOMBAR to copy the **gates**.

### 5.8. MULTIPLIER DECLARATIONS AND DEFINITIONS

Multipliers appear only in **structure diagrams** and MDECLARATIONTABLE records information about the **partial structures** in which they occur, and also the substituents to which they apply.

As each multiplier is defined, the values for it are placed in MDEFINITIONTABLE, and only on completion of the processing of a GENSAL sentence does **procedure** RECORDMULTS actually transfer this information to the ECTR, in appropriate FREQUENCY fields in the top bars of **child gates**.

5.9. TIDYING THE ECTR

Before returning to the main part of GENPROG, the interpreter calls **procedure** TIDYINTREP, which DISPOSES of certain redundant parts of the ECTR: these are mainly **partial structures** and their associated **gates** that were set up as parts of the entries in RDEFINITIONTABLE, and which were only required for checking purposes during **procedure** INTERPRET. A few other linked lists used for housekeeping purposes are also DISPOSED by TIDYINTREP.

**Procedure** OUTINTREP is used to output a representation of the ECTR to a diagnostics file, if desired, but this is intended only for programmer checks on the working of INTERPRET.

Finally, control is passed back to GENPROG, where the ECTR can be used for fragment generation and other purposes.

## CHAPTER 6

### CONCLUSIONS AND SUGGESTIONS FOR FURTHER WORK

---

"Give us the tools, and we will finish the job"

Churchill

#### 6.1. DEVELOPMENT OF A PATENT DOCUMENTATION SYSTEM

The input language GENSAL is essentially a means of describing generic chemical structures; it is not a means of encoding patent specifications as such. Many patents contain several generic structures covering, for example, the various components of a mixture, or different intermediates in a reaction pathway, and a patent documentation system might require all these to be encoded separately in GENSAL.

Considerable development work would be needed to take GENSAL, the ECTR, the interpreter program, and the work described by Welford<sup>67</sup> to form a comprehensive online patent information system, and further research work is also still required to make such a system viable.

From the interpreter point of view, the most immediate task is obviously to extend the program as it exists to implement the full GENSAL language, including conditions. But more important than this is likely to be the development of fragment-based, and possibly atom-by-atom search algorithms, and approaches to this are discussed by Welford.<sup>67</sup>

In order to evaluate such algorithms, a database of at least some hundreds of generic structures from patents, and some sample queries, will be required. A number of companies in the chemical and pharmaceutical industries have expressed interest in participating in the encoding of structures in GENSAL for this purpose.

The building up of a database will also enable a full evaluation to be made of the power of GENSAL to encode generic structures in patents. Additionally, it will permit an analysis of the effort required to encode a generic structure from a patent in GENSAL; Pötscher<sup>44</sup> has pointed out that in the encoding of generic structures in the GREMAS system, the difficult part of the operation is the analysis of the structure as described, rather than the selection of GREMAS terms, and this analysis would also

- to a large extent at least - need to be carried out for GENSAL coding.

Certainly, GENSAL coding from patents or abstracts is not a clerical task, though experience and a basic knowledge of elementary chemistry would be adequate qualifications for a coder. GENSAL coding is likely to require much less training than that required for encoding in a fragment-based system.

The possibility of automatic generation of GENSAL from patent specifications or abstracts is an interesting one; Nishida and Takamatsu<sup>193</sup> have recently described a method for extracting information from patent claim text, though their work was not related to chemical patents. The problem would be likely to be a very difficult one in the application of artificial intelligence techniques, and any system developed would certainly require human interference at points where the specification is ambiguous. Such automatic input of generic structures might however be essential if a viable back file of patents were to be built up.

Associated with the input of a large number of structures in GENSAL will be the need to add terms to the dictionary of nomenclatural terms, and this will require many decisions to be made as to the meanings to be assigned to vague terms, as discussed in Section 5.7.3.

## 6.2. OTHER POTENTIAL APPLICATIONS OF GENSAL

Whilst GENSAL has been designed for the encoding of the generic structures in patents, and thus to form part of an integrated patent information system, it has a number of potential applications outside the field of computer-based patent documentation systems, and these will now be mentioned briefly.

### 6.2.1. Non-Computer Description of Generic Structures

Since GENSAL is designed to be a complete and unambiguous means of describing generic chemical structures, it could well have a use in non-computer contexts, just as high-level programming languages such as Algol are often used for the non-computer description of algorithms.

GENSAL is intended to be readily comprehensible to a chemist or patent agent who has had a fairly minimal training in its use (though rather more training would be required to achieve efficiency and accuracy in encoding structures), and it might therefore have applications in printed abstracts of patents, or in current awareness bulletins. If such a printed publication were produced by a computer-typesetting process, then the use of GENSAL would give the added advantage of leaving a complete and unambiguous description of the generic structure in machine-readable form, so that it could, perhaps, be incorporated into a computerised storage and retrieval system at a later date.

The clarity and lack of ambiguity of GENSAL would make such descriptions much easier to understand than those currently found in patent specifications and abstracts.

It is even possible to speculate that GENSAL might ultimately be used for generic structure descriptions in the patent documents themselves, though this is likely to remain speculation for some time to come.

#### 6.2.2. Generic Structures in the Journal Literature

Figure 1.1 illustrated a generic structure from the journal literature, and such series of related compounds could conveniently be described using a single GENSAL structure, which might, if desired, be used for automatic generation of all the specific compounds covered, so that these could be registered in an appropriate specific-compound registry system. Integration with a quantitative structure activity relationship system might also allow the automatic identification of the compounds likely to be most active.

#### 6.2.3. Chemical Reaction Documentation

One of the problems in the documentation of chemical reactions is the description of the "generalised" reaction process. Normally this is done in terms of substructures for the reactant and the

product, which represent the "reaction centre" - i.e. the atoms and bonds actually involved in the reaction.

However, frequently the reaction is strongly influenced by the presence or absence of surrounding groups which do not actually participate in the bond changes. The description of the reaction centre as a generic structure, using GENSAL, would allow these variable surrounding groups to be taken into account, though the feasibility and development of a reaction indexing system based on this principle would need substantial research investigation.

#### 6.2.4. Specific Structure Search Queries

Many of the chemical structure search systems currently available commercially have some features for the use of generic structure queries in searches of files of specific structures. For the most part these allow only a very restricted type of generic structure, usually the specification of a few alternative atoms or groups at particular defined points in the query structure, though the COUSIN system at Upjohn,<sup>64</sup> described in Section 1.4.7, allows a greater degree of sophistication with its "R<sub>k</sub>" notation. Système DARC (Section 1.4.6) is also believed to be about to introduce substantial facilities for generic structure queries.

The use of GENSAL would permit much more complex generic structures to be input as queries for searches of a file of

specific structures, potentially without any need for modification of the search software. A GENSAL interpreter program would convert the GENSAL input to the ECTR internal representation, and from this a special fragment-generation module would produce a set of search fragments compatible with those normally used for searching the file, with appropriate "AND" and "OR" logic.

### 6.3. CONCLUSIONS

The work described in this Thesis forms a viable basis for an improved storage and retrieval system for generic structures in patents, and it is the hope of the author that it may be used in the development of such a system.

It is possible that, as discussed in this Chapter, the work may have applications in other areas also. Improved patent documentation systems may additionally have an effect on the processes of drafting and granting patents. In 1966 Frome<sup>194</sup> discussed the legal problems that could be caused by computer programs able to print out all the specific compounds covered by a generic structure. Blick<sup>195</sup> has also pointed out a similar problem with computer-aided synthesis packages, which could affect the patentability of synthesis routes suggested by such packages.

Whatever the fate of the present work, it is certain that storage and retrieval systems for generic chemical structures will have increasing importance in many areas for many years to come.

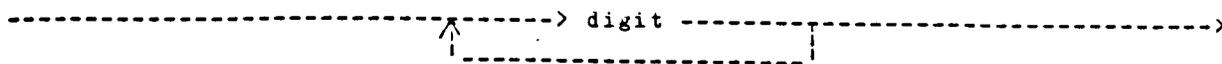
## APPENDIX 1

### GENSAL SYNTAX DIAGRAMS

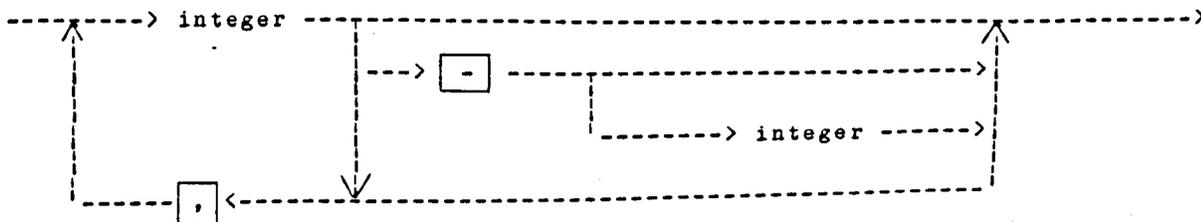
---

In these syntax diagrams, delimiter words and symbols are shown enclosed in boxes, data items are shown in upper case letters, and references to other syntax diagrams are shown in lower-case letters. A detailed discussion of the formal grammar of GENSAL is given in Section 3.13 of the text of this Thesis.

#### 1. integer



#### 2. integer range



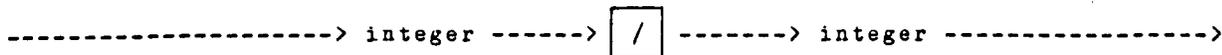
3. substituent



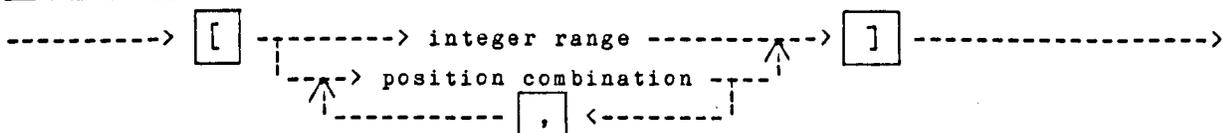
4. multiplier



5. position combination



6. position set



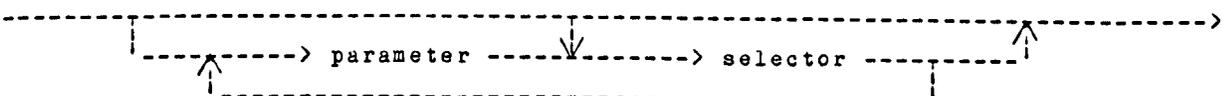
7. parameter



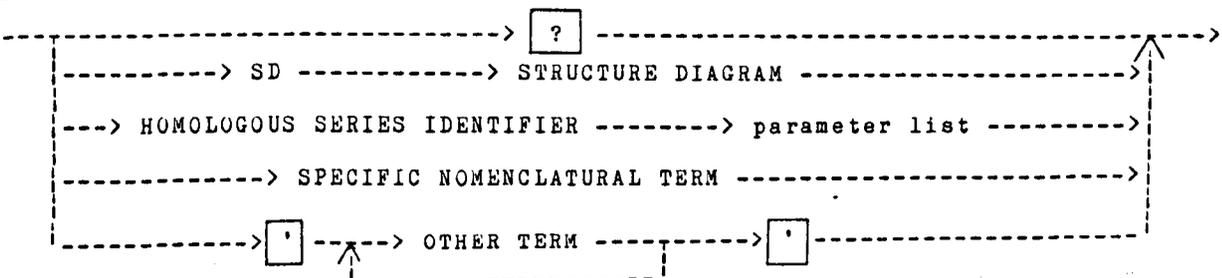
2. selector



9. parameter list



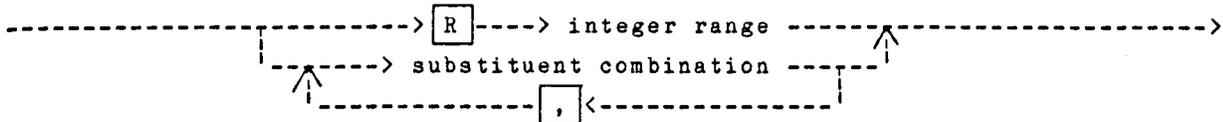
10. substituent value



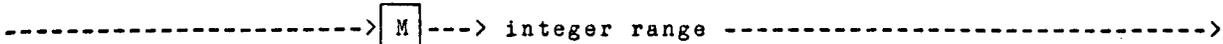
11. substituent combination



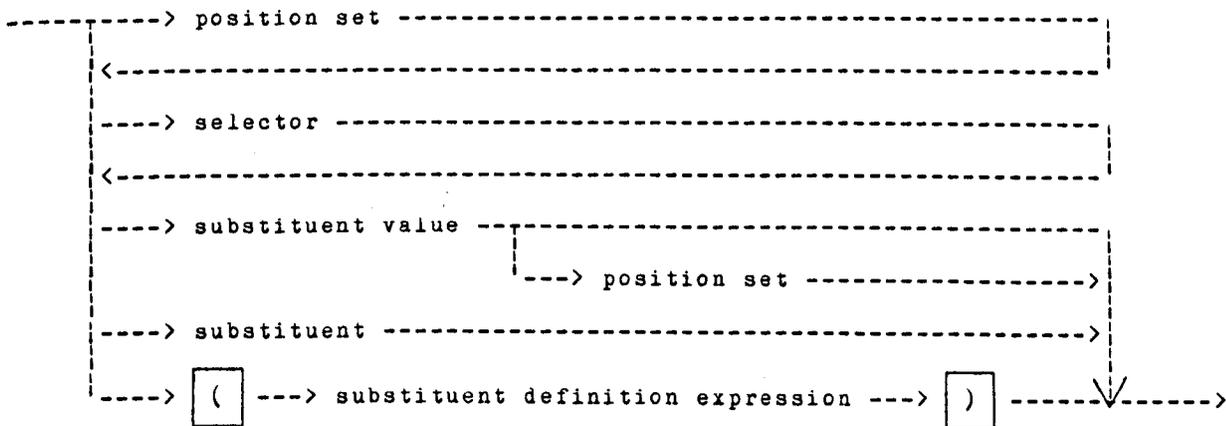
12. substituent group



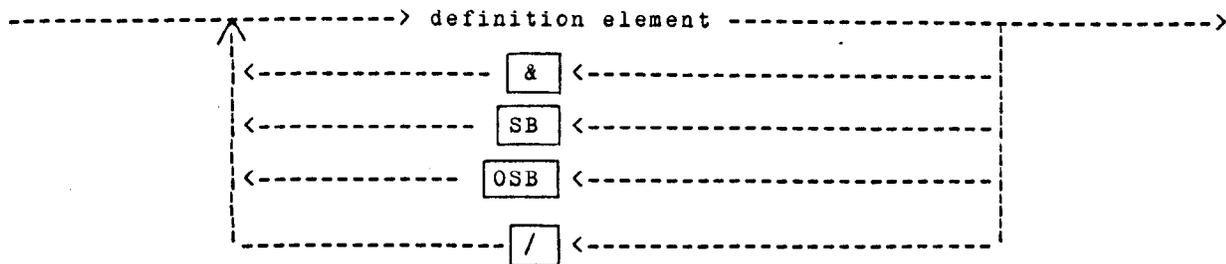
13. multiplier group



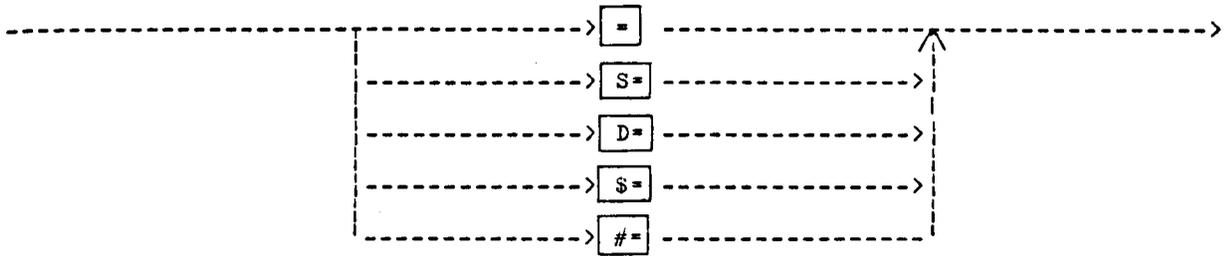
14. definition element



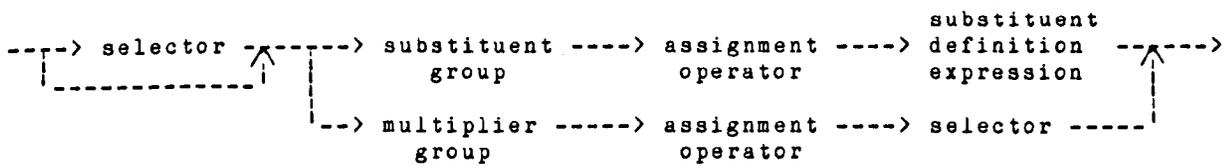
15. substituent definition expression



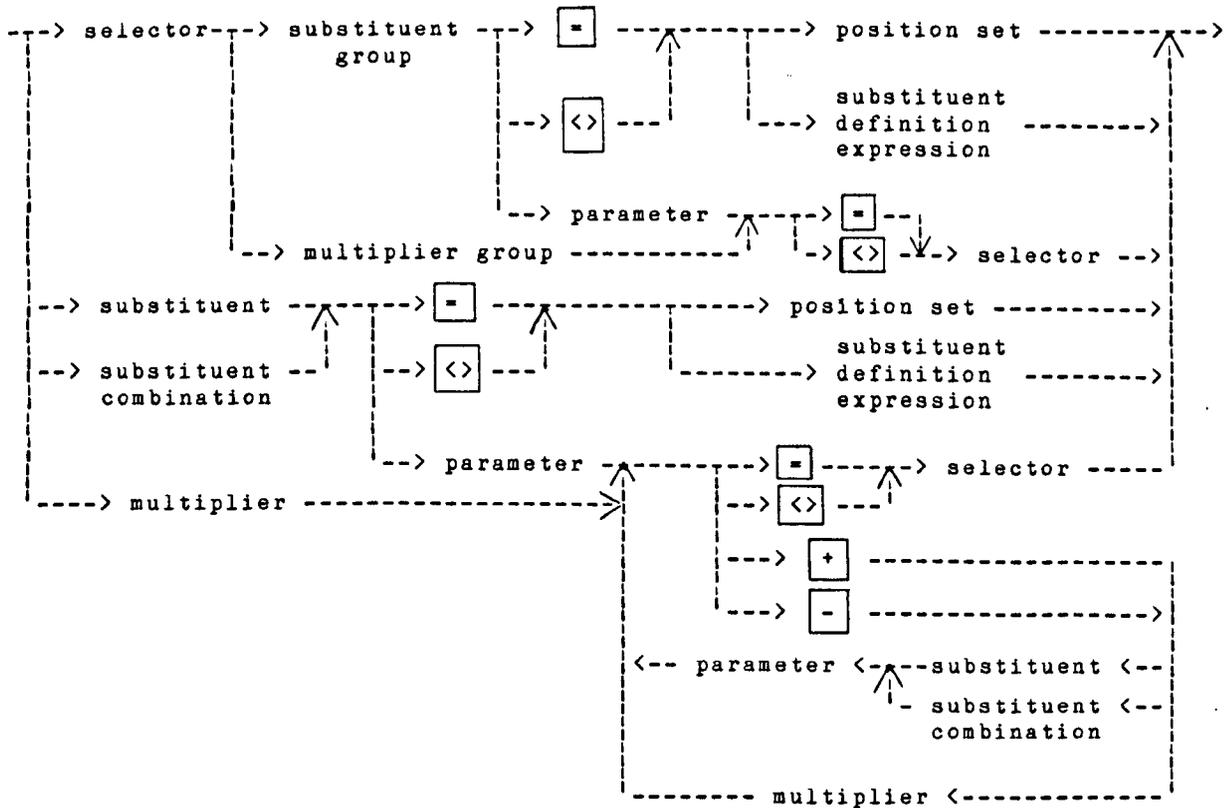
16. assignment operator



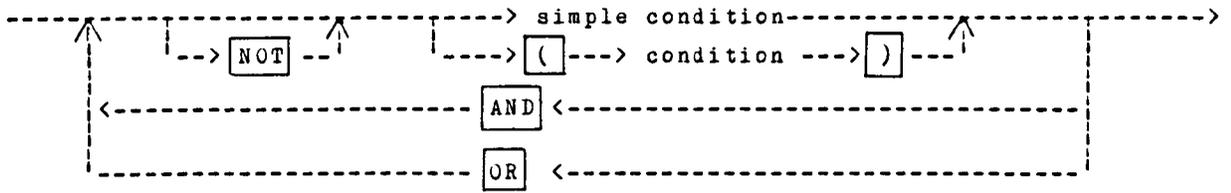
17. assignment statement



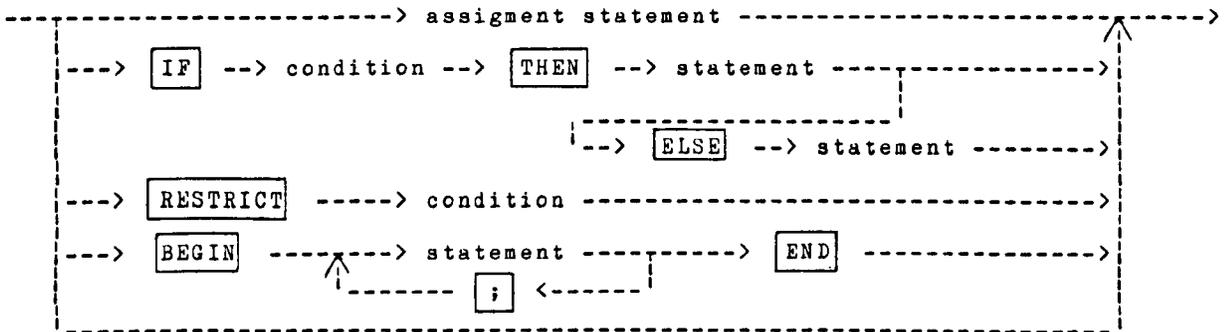
18. simple condition



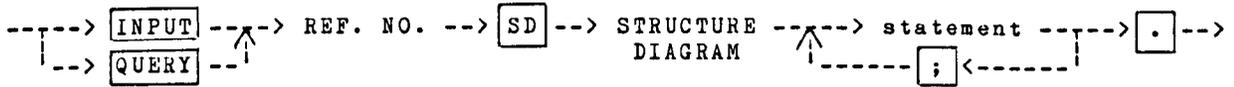
19. condition



20. statement



21. structure description



## APPENDIX 2

### BNF PRODUCTION RULES FOR GENSAL

This Appendix shows the Grammar of GENSAL using Backus-Naur Form (BNF) production rules. A slight variant of the "Extended BNF" metalanguage proposed by Wirth <sup>101</sup> is used, in which the syntactic constructs (non-terminal symbols) are shown enclosed in angle brackets, the symbol "::=" means "is replaced by", the symbol "|" means "or", curly brackets enclose symbols to be repeated zero or more times and square brackets enclose optional symbols. Terminal symbols included exactly as they stand are shown in **bold type** and are enclosed in double quote marks.

1. <arithmetic operator> ::= "+" | "-"
2. <assignment operator> ::= "=" | "S=" | "D=" | "\$=" | "#="
3. <assignment statement> ::= [ <selector> ]  
    <unselected assignment statement>
4. <character> ::= <letter> | <digit> | <special character>
5. <compound statement> ::= "BEGIN" <statement>  
    { ";" <statement> } "END"
6. <condition factor> ::= <simple condition> |  
    "(" <condition> ")" | "NOT" <condition factor>

7. <condition term> ::= <condition factor> { "AND"  
    <condition factor> }
8. <condition> ::= <condition term> { "OR" <condition term> }
9. <definition alternative> ::= <element combination> { <further  
    substitution operator> <element combination> }
10. <definition element> ::= [ <position set> ] [ <selector> ]  
    <unmodified definition element>
11. <definition relation> ::= <substituent variable>  
    <relational operator> <position set> |  
    <substituent variable> <relational operator>  
    <substituent definition expression>
12. <digit> ::= "0" | "1" | "2" | "3" | "4" | "5" | "6" | "7" |  
    "8" | "9"
13. <element combination> ::= <definition element> { "&"  
    <definition element> }
14. <empty> ::=
15. <further substitution operator> ::= "SB" | "OSB"
16. <group definition relation> ::= <substituent group>  
    <relational operator> <position set> |  
    <substituent group> <relational operator>  
    <substituent definition>
17. <group parameter relation> ::= <substituent group>  
    <parameter> <relational operator> <selector>
18. <group relation> ::= <substituent group relation> |  
    <multiplier group relation>
19. <homologous series identifier> ::= <nomenclature>
20. <IF statement> ::= "IF" <condition> "THEN" <statement>  
    [ "ELSE" <statement> ]
21. <initial character> ::= <letter> | <digit>
22. <integer range> ::= { <subrange> "," } <top range>
23. <integer term> ::= <multiplier> | <substituent variable>  
    <parameter>
24. <integer> ::= <digit> { <digit> }
25. <letter> ::= "A" | "B" | "C" | "D" | "E" | "F" | "G" | "H" |  
    "I" | "J" | "K" | "L" | "M" | "N" | "O" | "P" | "Q" |  
    "R" | "S" | "T" | "U" | "V" | "W" | "X" | "Y" | "Z"

26. <multiplier assignment> ::= <multiplier group>  
           <assignment operator> <selector>
27. <multiplier group relation> ::= <multiplier group>  
           <relational operator> <selector>
28. <multiplier group> ::= "M" <integer range>
29. <multiplier relation> ::= <multiplier>  
           { <arithmetic operator> <integer term> }  
           <relational operator> <selector>
30. <multiplier> ::= "M" <integer>
31. <nomenclature> ::= <initial character> { <character> }
32. <other term> ::= <nomenclature>
33. <parameter identifier> ::= "C" | "T" | "Q" | "E" | "Y" | "RC"  
           | "RN" | "RS" | "RA" | "Z"
34. <parameter list> ::= [ <selector> ]  
           { <parameter> <selector> }
35. <parameter relation> ::= <substituent variable> <parameter>  
           { <arithmetic operator> <integer term> }  
           <relational operator> <selector>
36. <parameter> ::= <parameter identifier> |  
           "" <substituent> ""
37. <position combination> ::= <integer> "/" <integer>
38. <position set> ::= "[" <positions> "]"
39. <positions> ::= <integer range> | <position combination>  
           { ", " <position combination> }
40. <reference number> ::= <integer>
41. <relational operator> ::= "=" | "<"
42. <restrict statement> ::= "RESTRICT" <condition>
43. <selector> ::= "<" <integer range> ">"
44. <simple condition> ::= <selector> <group relation> |  
           <substituent relation> | <multiplier relation>
45. <special character> (Implementation Dependent)
46. <specific nomenclatural term> ::= <nomenclature>
47. <statement> ::= <assignment statement> | <if statement> |  
           <restrict statement> | <compound statement> | <empty>

48. <structure description> ::= <structure type>  
     <reference number> "SD" <structure diagram> <statement>  
     { ";" <statement> } "."
49. <structure diagram> (Implementation Dependent)
50. <structure type> ::= "INPUT" | "QUERY"
51. <subrange> ::= <integer> | <integer> "-" <integer>
52. <substituent assignment> ::= <substituent group>  
     <assignment operator>  
     <substituent definition expression>
53. <substituent combination> ::= <substituent> "+" <substituent>
54. <substituent definition expression> ::=  
     <definition alternative>  
     { "/" <definition alternative> }
55. <substituent group relation> ::= <group definition relation>  
     <group parameter relation>
56. <substituent group> ::= "R" <integer range> |  
     <substituent combination>  
     { "," <substituent combination> }
57. <substituent relation> ::= <definition relation> |  
     <parameter relation>
58. <substituent value> ::= "?" | "SD" <structure diagram> |  
     <homologous series identifier> <parameter list > |  
     <specific nomenclatural term> | "" <other term>  
     { <other term> } ""
59. <substituent variable> ::= <substituent> |  
     <substituent combination>
60. <substituent> ::= "R" <integer>
61. <top range> ::= <subrange> | <integer> "-"
62. <unmodified definition element> ::= <substituent value> |  
     <substituent value> <position set> | <substituent> |  
     "(" <substituent definition expression> ")"
63. <unselected assignment statement> ::=  
     <substituent assignment> | <multiplier assignment>

**APPENDIX 3**

**GENSAL INTERPRETER PROGRAM**

```
1
2
3
4
5
6
7
8
9  PROCEDURE PROGERROR(ERRORCODE : INTEGER);
10 EXTERN;
11
12
13  PROCEDURE GOTOCOMMAND;
14 EXTERN;
15 { Sends the user to GENESIS command level via a GOTO in the main program}
16
17
18
19  PROCEDURE DESTROY(VAR PTR1 : PDOUBLIST);
20 EXTERN;
21 { This destroys the elements of a linked list of type PDOUBLIST, starting
22   at the element pointed to be the parameter PTR1, returned as NIL.
23   Called by INTERPRET\INTSET
24     INTERPRET\GROUPRANGE
25     INTERPRET\SELECTOR
26     INTERPRET\POSITIONSET
27     INTERPRET\MODIFYPOSITIONS\TRACEDOWNGATE
28     INTERPRET\ALTNVLIST\ELEMENT\SETCOMBARS\CHECKCOMBPOSNS
29     INTERPRET\ALTNVLIST\ELEMENT\PARAMETERLIST
30     INTERPRET\ASSIGNMENTSTMNT\MULTASSIGNMENT}
31
32
33
34
35  PROCEDURE ADDINTS (VAR PTR1      : PDOUBLIST;
36                    LOWER, UPPER : INTEGER);
37 EXTERN;
38 { Adds LOWER and UPPER to the values already in PTR1 (if they are contiguous),
39   or places them in a new DOUBLIST element, returned as PTR1, with the original
40   PTR1 as its NEXT field.
```

```
41      Called by INTERPRET\INTEGERRANGE\RANGEFRAGMENT
42      INTERPRET\SETINTS}
43
44
45
46  PROCEDURE PRINTNOM(NOMENVAL : STRING32);
47  EXTERN;
48  { Prints a nomenclatural term up to the last non-space character
49    Called by INTERPRET\ALTNVLIST\RECORDHELD}
50
51
52
53  PROCEDURE DELETEDGENSAL(VAR LINE1 : PLINELIST);
54  EXTERN;
55  { Deletes a linked list of GENSAL lines, headed by LINE1, which is returned
56    with value NIL.
57    Called by INTERPRET\ALTNVLIST\ELEMENT\TRANSLATENOMEN}
58
59
60
61  PROCEDURE DECODECT (VAR CTLINE : PLINELIST;
62                    DISPLAYING : BOOLEAN);
63  EXTERN;
64  { Decodes a connection table from character-string format, beginning in CTLINE^.
65    LINE, making entries in FELDCT and FELDBD. CTLINE is left pointing at the last
66    line of the connection table string. FELDMN is used to display the structure
67    diagram if DISPLAYING.
68    Called by INTERPRET\READSD}
69
70
71
72  PROCEDURE ENCODECT(VAR CTLINE : PLINELIST);
73  EXTERN;
74  { Encodes the contents of FELDCT and FELDBD as a character string, and
75    places it in successive lines, starting with CTLINE, which is
76    returned pointing to the last line of encoded connection table.
77    Called by INTERPRET\READSD}
78
79
80
```

```

81 FUNCTION NORECORD(NOMEN      : STRING32;
82                   VAR ADDRESS : INTEGER): BOOLEAN;
83 EXTERN;
84 { Checks whether or not there is a record held in SPSDICT for NOMEN,
85   The ADDRESS from SPSDICT for the term is returned as a side effect.
86   Called by INTERPRET\ALTNVLIST\RECORDHELD}
87
88
89
90 FUNCTION TERMREAD(VAR TERM : STRING32) : BOOLEAN;
91 EXTERN;
92 { Reads a single TERM from the terminal, upper-cases it, and returns FALSE
93   if it has no characters.
94   Called by INTERPRET\ALTNVLIST\RECORDHELD}
95
96
97
98
99 PROCEDURE LISTPARAMS(VAR OUTFILE : TEXT;
100                   PARAMLIST  : TPARAMLIST);
101 EXTERN;
102 { Lists the parameters in PARAMLIST in file OUTFILE, which must
103   already have been RESET.
104   Called by INTERPRET\OUTINTREP\WRITEPS}
105
106
107
108 FUNCTION SPSVARIETY(ADDRESS      : INTEGER;
109                   DISPLAYING    : BOOLEAN) : TPSVARIETY;
110 EXTERN;
111 { Returns the variety of partial structure, whose record begins at ADDRESS in
112   SPSFILE, optionally DISPLAYING the structure. SPECIFIC PSs are entered in
113   FELDCT/FELDBD, GENERIC PSs in SPSPARAMLIST and OTHER PSs in INSERTGENEX.
114   Called by INTERPRET\ALTNVLIST\ELEMENT\TRANSLATENOMEN}
115
116
117
118 PROCEDURE READFELDMANN;
119 EXTERN;
120 { Reads the Feldmann table from FELDFIL.

```

```

121      Called by INTERPRET\PROCESSCT}
122
123
124
125  PROCEDURE INTERPRET(VAR FIRSTLINE : PLINELIST;
126                      VAR ECTRSIZE  : INTEGER;
127                      INTERACTIVE   : BOOLEAN);
128
129  { This is the GENSAL interpreter routine, and performs syntactic and semantic
130    analysis on a GENSAL sentence, creating the ECTR. }
131
132  CONST NOTFIXED      = 0;
133
134  TYPE DELIMTYPE      =(INVALIDTOKEN,GAMPERSAND,GPRIME,GLPAREN,GRPAREN,GPLUS,GCOMMA,
135                      GHYPHEN,GPERIOD,GSLASH,GSEMI,GOPENANG,GNOTEQ,GEQUALS,
136                      GCLOSANG,GQUEST,GLSQUARE,GRSQUARE,GAND,GBEGIN,GC,GE,GELSE,
137                      GEND,GIF,GINPUT,GM,GN,GOR,GORBY,GOSB,GP,GQ,GQUERY,GR,GRA,
138                      GRC,GRESTRICT,GRF,GRN,GRS,GSB,GSD,GT,GTHEN,GY,GZ,GDEQ,GSEQ,
139                      GHASHEQ,GDOLEQ);
140  TOKENNATURE         =(DELIMITER,INTEGRAL,NOMENCLATURE);
141  TINPUTMODE          =(TERMINAL, STOREDGENSAL, INSERTTEXT);
142  TBOND MAG           = 0..3;
143  DELIMSET            = SET OF DELIMTYPE;
144  TOKENTYPE           = RECORD
145                      CASE NATURE : TOKENNATURE OF
146                          DELIMITER   : (DELIMVAL : DELIMTYPE);
147                          INTEGRAL    : (INTEGVAL : INTEGER);
148                          NOMENCLATURE : (NOMENVAL : STRING32)
149                      END;
150  PTOKENLIST          = ^TTOKENLIST;
151  TTOKENLIST          = RECORD
152                      TOKENVAL : TOKENTYPE;
153                      NEXT     : PTOKENLIST
154                      END;
155  PPSLIST              = ^PSLIST;
156  PSLIST              = RECORD
157                      PARSTRUCT,
158                      FURTHERSUB   : PTRPSTYPE;
159                      COMBINS      : PCOMBINLIST;
160                      CONNBONDS    : TCONNBONDS;

```

```

161             PRNTPOSNS : PTGROUPMEMS;
162             COPYCHILDPS : BOOLEAN;
163             NEXT       : PPSLIST
164             END;
165     PMDECLIST = ^MDECLIST;
166     MDECLIST  = RECORD
167             SUBSTDECN : PPSLIST;
168             NEXT      : PMDECLIST
169             END;
170
171     VAR TOKEN           : TOKENTYPE;
172     TOKENLIST          : PTOKENLIST;
173     LINENUMBER         : INTEGER;
174     VALIDLENGTH       : 0..MAXLENGTH;
175     DEFNMULT,
176     DEFNSUBS,          { substituents so far defined }
177     DECLMULT,
178     DECLSUBS          : INTEGSET;
179     CONDITIONSPRESENT : BOOLEAN;
180     RDECLARATIONTABLE : ARRAY[SUBSTITUENT] OF PPSLIST;
181     RDEFINITIONTABLE  : ARRAY[SUBSTITUENT] OF PCOMBINLIST;
182     MDECLARATIONTABLE : ARRAY[MULTIPLIER] OF PMDECLIST;
183     MDEFINITIONTABLE  : ARRAY[MULTIPLIER] OF INTRECORD;
184     CURRENTLINE       : PLINELIST;
185     SUBST              : SUBSTITUENT;
186     ZEROFREQ,
187     ESSENTFREQ,
188     OPTFREQ           : PDOUBLIST;
189     BONDMATCHARRAY    : ARRAY[BONDORDER] OF PACKED ARRAY[BONDORDER] OF BONDORDER;
190     BONDSTRING        : ARRAY[BONDORDER] OF STRING2;
191     INPUTMODE         : TINPUTMODE;
192     INSERTHSTPS       : PTRPSTYPE;
193     IRLISTBOT         : PIRLIST;
194
195
196
197     {-----}
198     PROCEDURE INITIALISE;
199
200     { Sets initial values for variables.

```

```

201         Called by Body of INTERPRET}
202
203     VAR SUBST      : SUBSTITUENT;
204         BOND       : BONDORDER;
205         MULT       : MULTIPLIER;
206         BONDFILE   : FILE OF PACKED ARRAY[BONDORDER] OF BONDORDER;
207
208
209
210     FUNCTION NEWFREQ(ONE, TWO : INTEGER) : PDOUBLIST;
211
212     VAR NF : PDOUBLIST;
213
214     BEGIN
215     NEW(NF);
216     NF^.FIRST := ONE;
217     NF^.SECOND := TWO;
218     NF^.NEXT := NIL;
219     NEWFREQ := NF
220     END;
221
222
223
224     BEGIN {Body of INITIALISE}
225     ESSENTFREQ := NEWFREQ(1,1);
226     OPTFREQ := NEWFREQ(0,1);
227     ZEROFREQ := NEWFREQ(0,0);
228     ECTRSIZE := 18;
229     CONDITIONSPRESENT := FALSE;
230     WRITELN;
231     WRITELN;
232     LINENUMBER := 0;
233     DECLSUBS := [];
234     DECLMULT := [];
235     DEFNMULT := [];
236     DEFNSUBS := [];
237     IF INTERACTIVE THEN INPUTMODE := TERMINAL
238         ELSE INPUTMODE := STOREDGENSAL;
239     INTERNALREP.CONSTANTPART := NIL;
240     TOKENLIST := NIL;

```

```

241 N := MAXLENGTH;
242 FOR SUBST := 1 TO MAXVARS DO
243 BEGIN
244 RDECLARATIONTABLE[SUBST] := NIL;
245 RDEFINITIONTABLE[SUBST] := NIL
246 END;
247 FOR MULT := 1 TO MAXVARS DO
248 BEGIN
249 MDECLARATIONTABLE[MULT] := NIL;
250 MDEFINITIONTABLE[MULT].TOPRANGE := NOTSET;
251 MDEFINITIONTABLE[MULT].SUBRANGES := NIL
252 END;
253 CURRENTLINE := NIL;
254
255 RESET(BONDFILE, 'LI2GEN>BONDFILE');
256 FOR BOND := NOTSPECIFIED TO RINGTAUT DO READ(BONDFILE, BONDMATCHARRAY[BOND]);
257 RESET(BONDFILE, '@TTY');
258
259 BONDSTRING[NOTSPECIFIED] := 'NS';
260 BONDSTRING[ANY] := 'A';
261 BONDSTRING[CHAIN] := 'C';
262 BONDSTRING[RING] := 'R';
263 BONDSTRING[SINGLE] := 'S';
264 BONDSTRING[DOUBLE] := 'D';
265 BONDSTRING[TRIPLE] := 'T';
266 BONDSTRING[CHASING] := 'CS';
267 BONDSTRING[CHAI DOUB] := 'CD';
268 BONDSTRING[CHAITRIP] := 'CT';
269 BONDSTRING[CHAITAUT] := 'TC';
270 BONDSTRING[RINGSING] := 'RS';
271 BONDSTRING[RING DOUB] := 'RD';
272 BONDSTRING[RINGTRIP] := 'RT';
273 BONDSTRING[AROMATIC] := 'RA';
274 BONDSTRING[RINGTAUT] := 'TR'
275
276 END; { of INITIALISE
277 -----}
278
279
280

```

```

281 PROCEDURE WRITEMESSAGE(ERRORCODE,
282                          NUMDATA   : INTEGER;
283                          STRINGDATA : STRING4);
284
285 { Obtains an error message from LI2GEN>ERRORMSGs, and prints it at the
286   terminal, interposing data where necessary.
287   Called by FAILURE
288         ERROR
289         PROCESSCT\REJECT}
290
291 VAR STRINGPOS : 1..5;
292     MSGCHAR   : CHAR;
293     LINE      : INTEGER;
294
295 BEGIN
296   RESET(INPUT, 'LI2GEN>ERRORMSGs');
297   FOR LINE := 1 TO (ERRORCODE-1) DO READLN;
298   STRINGPOS := 1;
299   WHILE NOT EOLN(INPUT) DO
300     BEGIN
301       READ(MSGCHAR);
302       CASE MSGCHAR OF
303         '#'   : WRITE(NUMDATA : 1);
304         '$'   : BEGIN
305                   WRITE(STRINGDATA[STRINGPOS]);
306                   STRINGPOS := STRINGPOS + 1
307                 END;
308         OTHERWISE WRITE(MSGCHAR)
309       END
310     END;
311   WRITELN;
312   RESET(INPUT, '@TTY')
313 END;
314
315
316
317 PROCEDURE FAILURE(ERRORCODE,
318                  NUMDATA   : INTEGER;
319                  STRINGDATA : STRING4);
320

```

```

321 { Called when an irrecoverable error is encountered, and processing cannot
322 continue. A message is printed, and the use retruned to GENESIS command
323 mode. }
324
325 BEGIN
326 WRITELN;
327 WRITELN('**** FAILURE ', ERRORCODE : 2);
328 WRITEMESSAGE(ERRORCODE, NUMDATA, STRINGDATA);
329 WRITELN;
330 WRITELN('Edit existing GENSAL or start again!');
331 WRITELN;
332 GOTOCOMMAND
333 END;
334
335
336
337 PROCEDURE REDUCECTR(PTR : PDOUBLIST);
338
339 BEGIN
340 WHILE PTR <> NIL DO
341 BEGIN
342 ECTRSIZE := ECTRSIZE - 6;
343 PTR := PTR^.NEXT
344 END
345 END;
346
347
348
349 {-----}
350
351 PROCEDURE GETTOKEN
352 THE LEXICAL ANALYZER
353 -----}
354
355
356
357 PROCEDURE GETTOKEN;
358
359 { Places the next token in the GENSAL input stream in TOKEN.
360 Called by NEXTTOKEN

```

```

361             LOOKAHEAD
362             ERROR}
363
364 VAR M         : 1..MAXLENGTH;
365     STARTED   : BOOLEAN;
366
367
368
369 PROCEDURE READLINE;
370
371 { Reads one line of GENSAL input according to INPUTMODE, checking for TERMINAL
372   that it contains no more than 99 characters, and building up the linked list
373   of lines. For all INPUTMODE replaces all lower-case alphabetic by upper-case. }
374
375 LABEL 10;
376
377 VAR CH        : CHAR;
378     M         : 0..MAXLENGTH;
379     NEWLINE    : PLINELIST;
380
381 BEGIN
382 CASE INPUTMODE OF
383     TERMINAL   : BEGIN
384                 LINENUMBER := LINENUMBER + 1;
385                 10 : WRITE(LINENUMBER : 3);
386                 IF DIAGNOSTICS THEN WRITE(ECTRSIZE : 7 );
387                 WRITE(' GENSAL: ');
388                 READLN(BUFFER : N);
389                 IF N=0 THEN
390                     BEGIN
391                         WRITELN;
392                         WRITELN('GENSAL input terminated by user.');
```

```

401         FOR M := (MAXLENGTH-12) TO (MAXLENGTH-1) DO WRITE(BUFFER[M]);
402         WRITELN('".');
403         REPEAT
404             WRITE('                                OK? (Y/N) > ');
405             READLN(CH)
406         UNTIL (CH='Y') OR (CH='y') OR (CH='N') OR (CH='n');
407         M := MAXLENGTH;
408         REPEAT
409             BUFFER[M] := ' ';
410             M := M-1
411         UNTIL (BUFFER[M]=' ') OR (M=0) OR (CH='Y') OR (CH='y');
412         IF (CH='N') OR (CH='n') THEN WRITELN('Line truncated from last space.')
```

APPENDIX 3:

```

413         END;
414     NEW(NEWLINE);
415     WITH NEWLINE^ DO
416     BEGIN
417         LAST := CURRENTLINE;
418         NEXT := NIL;
419         LINE := BUFFER
420     END;
421     IF CURRENTLINE = NIL
422     THEN FIRSTLINE := NEWLINE
423     ELSE CURRENTLINE^.NEXT := NEWLINE;
424     CURRENTLINE := NEWLINE
425 END;
426
427 STOREDGENSAL : BEGIN
428     LINENUMBER := LINENUMBER + 1;
429     IF LINENUMBER = 1
430     THEN NEWLINE := FIRSTLINE
431     ELSE NEWLINE := CURRENTLINE^.NEXT;
432     IF NEWLINE = NIL
433     THEN BEGIN
434         WRITELN;
435         WRITELN('End of stored GENSAL.');
```

GENSAL INTERPRETER

```

436         WRITELN('Input at the terminal:');
437         WRITELN;
438         INPUTMODE := TERMINAL;
439         GOTO 10
440     END
```

```

441         ELSE CURRENTLINE := NEWLINE;
442         WRITE(LINENUMBER : 3);
443         IF DIAGNOSTICS THEN WRITE(ECTRSIZE : 7);
444         WRITE(' GENSAL: ');
445         BUFFER := CURRENTLINE^.LINE;
446         WHILE (BUFFER[N]=' ') AND (N>1) DO N := N-1;
447         IF (N=1) AND (BUFFER[1]=' ') THEN N:= 0;
448         FOR M := 1 TO N DO WRITE(BUFFER[M]);
449         WRITELN
450     END;
451
452     INSERTTEXT : BEGIN
453         IF CURRENTLINE = NIL THEN
454             PROGERROR(1); {Unterminated GENSAL expression in SPSfile}
455             BUFFER := CURRENTLINE^.LINE;
456             CURRENTLINE := CURRENTLINE^.NEXT
457         END
458     END;
459
460     FOR M:= 1 TO MAXLENGTH DO
461         IF (BUFFER[M] >= 'a') AND (BUFFER[M] <= 'z')
462             THEN BUFFER[M] := CHR( ORD(BUFFER[M]) - ORD('a') + ORD('A') );
463     N := 1
464     END (* Of READLINE *);
465
466
467
468     FUNCTION CHECK (TESTDELIM : DELIMTYPE) : BOOLEAN;
469
470     { Returns TRUE if the delimiter passed as TESTDELIM is found, correctly
471       terminated in BUFFER }
472
473     VAR RESULT      : (NOTFOUND, PENDING, FOUND);
474         M           : 0..MAXLENGTH;
475         TERMCHARS   : SET OF CHAR;
476         DELIMSTRING : ALFA;
477
478     BEGIN
479     CASE TESTDELIM OF
480         GAND:      DELIMSTRING := 'AND      ';

```

```

481      GBEGIN:      DELIMSTRING := 'BEGIN  '
482      GC:          DELIMSTRING := 'C      '
483      GE:          DELIMSTRING := 'E      '
484      GELSE:       DELIMSTRING := 'ELSE   '
485      GEND:        DELIMSTRING := 'END    '
486      GIF:         DELIMSTRING := 'IF     '
487      GINPUT:      DELIMSTRING := 'INPUT  '
488      GM:          DELIMSTRING := 'M      '
489      GN:          DELIMSTRING := 'N      '
490      GOR:         DELIMSTRING := 'OR     '
491      GORBY:       DELIMSTRING := 'ORBY   '
492      GOSB:        DELIMSTRING := 'OSB    '
493      GP:          DELIMSTRING := 'P      '
494      GQ:          DELIMSTRING := 'Q      '
495      GQUERY:      DELIMSTRING := 'QUERY  '
496      GR:          DELIMSTRING := 'R      '
497      GRA:         DELIMSTRING := 'RA     '
498      GRC:         DELIMSTRING := 'RC     '
499      GRESTRIC:    DELIMSTRING := 'RESTRICT'
500      GRF:         DELIMSTRING := 'RF     '
501      GRN:         DELIMSTRING := 'RN     '
502      GRS:         DELIMSTRING := 'RS     '
503      GSB:         DELIMSTRING := 'SB     '
504      GSD:         DELIMSTRING := 'SD     '
505      GT:          DELIMSTRING := 'T      '
506      GTHEN:       DELIMSTRING := 'THEN   '
507      GY:          DELIMSTRING := 'Y      '
508      GZ:          DELIMSTRING := 'Z      '
509      GDEQ:        DELIMSTRING := 'D=    '
510      GSEQ:        DELIMSTRING := 'S=    '
511      GHASHEQ:     DELIMSTRING := '#=    '
512      GDOLEQ:     DELIMSTRING := '$=    '
513      END;
514
515      TERMCHARS := [' ', '#', '$', '''..''), '. '..'9', '; '..'?', '[', ']' ];
516      M := 0;
517      RESULT := PENDING;
518      WHILE RESULT = PENDING DO
519          IF M=8
520              THEN IF BUFFER[N+M] IN TERMCHARS

```

```

521         THEN RESULT := FOUND
522         ELSE RESULT := NOTFOUND
523     ELSE IF (BUFFER[N+M] = DELIMSTRING [M+1])
524         THEN (* match found *) IF BUFFER[N+M] = ' '
525             THEN RESULT := FOUND (* i.e. match is on the space *)
526             ELSE M:=M+1 (* delimiter is still being read *)
527         ELSE (* no match *) IF(DELIMSTRING[M+1] <> ' ')
528             THEN RESULT := NOTFOUND (* not end of delimiter *)
529             ELSE IF BUFFER[N+M] IN TERMCHARS
530                 THEN RESULT := FOUND (* terminated *)
531                 ELSE IF TESTDELIM IN [GDEQ..GDOLEQ]
532                     THEN RESULT := FOUND (* no termination needed *)
533                     ELSE RESULT := NOTFOUND (* not terminated *);
534     IF RESULT = FOUND
535     THEN BEGIN
536         CHECK := TRUE;
537         N := N + M + ORD(M=7)
538     END
539     ELSE CHECK := FALSE
540 END;

541
542
543 PROCEDURE FINDNOMEN(VAR NOMENVAL : STRING32);
544
545 (* Extracts characters from BUFFER until nomenclature is correctly terminated.
546    If there are less than 32 characters before termination, then NOMENVAL
547    is packed with spaces; if more then the excess is discarded. A number of
548    of right parentheses equal to the number of left parentheses encountered
549    is accepted before a right parenthesis terminates the nomenclature. *)
550
551 VAR TERMINATED   : BOOLEAN;
552     M             : 1..32;
553     BRACKETCOUNT : 0..MAXLENGTH;
554     TERMCHARS    : SET OF CHAR;
555
556 BEGIN
557     TERMCHARS := [';', ',', ' ', '[', '/', '""', '<', '.'];
558     TERMINATED := FALSE;
559     BRACKETCOUNT := 0;
560     REPEAT

```

```

561     IF BUFFER[N] IN TERMCHARS
562     THEN TERMINATED := TRUE
563     ELSE IF BUFFER[N]='('
564         THEN BRACKETCOUNT := BRACKETCOUNT+1
565     ELSE IF BUFFER[N]=')'
566         THEN IF BRACKETCOUNT>0
567             THEN BRACKETCOUNT := BRACKETCOUNT-1
568             ELSE TERMINATED := TRUE;
569     IF NOT TERMINATED THEN N := N+1
570 UNTIL TERMINATED;
571 FOR M := 1 TO 32 DO
572     IF (M+VALIDLENGTH) < N
573     THEN NOMENVAL[M] := BUFFER[M+VALIDLENGTH]
574     ELSE NOMENVAL[M] := ' ';
575     TOKEN.NATURE := NOMENCLATURE
576 END;
577
578
579
580
581 FUNCTION CHECKINT : BOOLEAN;
582
583 { Returns TRUE if the token beginning at the current position in BUFFER is
584   an integer, and not nomenclature beginning with a digit. It checks this
585   by seeing if any leading digits, hyphens and commas are followed by an
586   alphabetic letter other than an R alone (as in substituent groups). }
587
588 VAR M          : 0..MAXLENGTH;
589
590 BEGIN
591     M := 0;
592     WHILE BUFFER[N+M] IN ['0'..'9', '-', ','] DO M := M+1;
593     IF BUFFER[N+M] IN ['A'..'Z']
594     THEN IF BUFFER[N+M] = 'R'
595         THEN CHECKINT := NOT (BUFFER[N+M+1] IN ['A'..'Z'])
596         ELSE CHECKINT := FALSE
597     ELSE CHECKINT := TRUE
598 END;
599
600

```

```

601     FUNCTION EXTRACTINT : INTEGER;
602
603     { Returns the integer at the current position in BUFFER }
604
605     VAR INTBUFF : ARRAY[1..9] OF CHAR;
606         INT,M,J,
607         K,MULT : INTEGER;
608
609     BEGIN
610     M := 0;
611     WHILE (BUFFER[N+M] IN ['0'..'9']) AND (M<9) DO
612     BEGIN
613     INTBUFF[M+1] := BUFFER [N+M];
614     M := M+1
615     END;
616     INT := 0;
617     FOR J := 0 TO (M-1) DO
618     BEGIN
619     MULT := 1;
620     FOR K := 1 TO J DO MULT := MULT*10;
621     INT := INT + MULT * ( ORD (INTBUFF[M-J]) -ORD('0'))
622     END;
623     EXTRACTINT := INT;
624     TOKEN.NATURE := INTEGRAL;
625     N := N + M
626     END;
627
628
629
630
631
632
633     BEGIN      (* Body of Procedure GETTOKEN *)
634
635     REPEAT
636     IF N=MAXLENGTH THEN READLINE;
637     WHILE (BUFFER[N]=' ') AND (N<MAXLENGTH) DO N := N+1;
638     STARTED := N<MAXLENGTH;
639     UNTIL STARTED;
640     VALIDLENGTH := N-1;

```

```

641 WITH TOKEN DO
642 IF BUFFER[N] IN ['A'..'E', 'I', 'M'..'T', 'Y', 'Z', '#', '$', '""..'')', '+..'/'', ';..'?'', '[', ']']
643 THEN BEGIN
644     CASE BUFFER[N] OF
645     'A': DELIMVAL := GAND;
646     'B': DELIMVAL := GBEGIN;
647     'C': DELIMVAL := GC;
648     'D': DELIMVAL := GDEQ;
649     'E': IF BUFFER[N+1] = 'L'
650         THEN DELIMVAL := GELSE
651         ELSE IF BUFFER[N+1] = 'N'
652             THEN DELIMVAL := GEND
653             ELSE DELIMVAL := GE;
654     'I': IF BUFFER[N+1] = 'F' THEN DELIMVAL := GIF
655         ELSE DELIMVAL := GINPUT;
656     'M': DELIMVAL := GM;
657     'N': DELIMVAL := GN;
658     'O': IF BUFFER[N+1] = 'R'
659         THEN IF BUFFER[N+2] = 'B' THEN DELIMVAL := GORBY
660             ELSE DELIMVAL := GOR
661         ELSE DELIMVAL := GOSB;
662     'P': DELIMVAL := GP;
663     'Q': IF BUFFER[N+1] = 'U' THEN DELIMVAL := GQUERY
664         ELSE DELIMVAL := GQ;
665     'R': IF BUFFER[N+1] IN ['A', 'C', 'E', 'N', 'S']
666         THEN CASE BUFFER[N+1] OF
667         'A' : DELIMVAL := GRA;
668         'C' : DELIMVAL := GRC;
669         'E' : DELIMVAL := GRESTRIC;
670         'F' : DELIMVAL := GRF;
671         'N' : DELIMVAL := GRN;
672         'S' : DELIMVAL := GRS
673         END
674         ELSE DELIMVAL := GR;
675     'S': IF BUFFER[N+1] = 'B'
676         THEN DELIMVAL := GSB
677         ELSE IF BUFFER[N+1] = 'D'
678             THEN DELIMVAL := GSD
679             ELSE DELIMVAL := GSEQ;
680     'T': IF BUFFER[N+1] = 'H' THEN DELIMVAL := GTHEN

```

APPENDIX 3:

GENSAL INTERPRETER

```

681                                     ELSE DELIMVAL := GT;
682 'Y': DELIMVAL := GY;
683 'Z': DELIMVAL := GZ;
684 '&': DELIMVAL := GAMPERSAND;
685 ''': DELIMVAL := GPRIME;
686 '(' : DELIMVAL := GLPAREN;
687 ')' : DELIMVAL := GRPAREN;
688 '+' : DELIMVAL := GPLUS;
689 ',' : DELIMVAL := GCOMMA;
690 '-' : DELIMVAL := GHYPHEN;
691 '.' : DELIMVAL := GPERIOD;
692 '/' : DELIMVAL := GSLASH;
693 ';' : DELIMVAL := GSEMI;
694 '<': IF BUFFER[N+1] = '>' THEN DELIMVAL := GNOTEQ
695                                     ELSE DELIMVAL := GOPENANG;
696 '=' : DELIMVAL := GEQUALS;
697 '>': DELIMVAL := GCLOSANG;
698 '?': DELIMVAL := GQUEST;
699 '[' : DELIMVAL := GLSQUARE;
700 ']' : DELIMVAL := GRSQUARE;
701 '#' : DELIMVAL := GHASHEQ;
702 '$': DELIMVAL := GDOLEQ
703 END (* of case *);
704 IF DELIMVAL >= GAND
705     THEN IF NOT CHECK(DELIMVAL) THEN FINDNOMEN(NOMENVAL)
706                                     ELSE NATURE := DELIMITER
707     ELSE BEGIN
708         NATURE := DELIMITER;
709         IF DELIMVAL=GNOTEQ THEN N := N + 2
710                                     ELSE N := N + 1
711     END
712 END (* of IF BUFFER[N] THEN *)
713 ELSE IF BUFFER[N] IN ['0'..'9']
714     THEN IF CHECKINT THEN INTEGVAL := EXTRACTINT
715                                     ELSE FINDNOMEN(NOMENVAL)
716     ELSE FINDNOMEN(NOMENVAL)
717 END;
718
719
720 { END OF PROCEDURE GETTOKEN (THE LEXICAL ANALYZER)

```

```
721 -----}
722
723
724
725
726 PROCEDURE NEXTTOKEN;
727
728 { Obtains the next token, either from the queue of tokens already produced
729   by LOOKAHEAD, or by a direct call to GETTOKEN. }
730
731 VAR TPTR : PTOKENLIST;
732
733 BEGIN
734   IF TOKENLIST = NIL
735     THEN GETTOKEN
736     ELSE BEGIN
737         TOKEN := TOKENLIST^.TOKENVAL;
738         TPTR := TOKENLIST;
739         TOKENLIST := TOKENLIST^.NEXT;
740         DISPOSE(TPTR)
741     END
742 END;
743
744
745
746 PROCEDURE LOOKAHEAD;
747
748 { If TOKENLIST is NIL (i.e. this is the first lookahead) then the current TOKEN
749   is placed at the bottom of a queue of tokens (TOKENLIST). GETTOKEN is used to
750   obtain a new token from the input stream, which is also added to the bottom
751   of the queue. The next call to NEXTTOKEN will therefore restore the original
752   TOKEN, and the subsequent call will return the following token.
753   Called by ALTNVLIST\POSITIONSET
754         ASSIGNMENTSTMNT\SUBSTGROUP}
755
756 VAR TOKENPTR : PTOKENLIST;
757
758 BEGIN
759   IF TOKENLIST=NIL
760     THEN BEGIN
```

```

761         NEW(TOKENLIST);
762         TOKENLIST^.TOKENVAL := TOKEN;
763         TOKENLIST^.NEXT := NIL
764     END;
765     TOKENPTR := TOKENLIST;
766     WHILE TOKENPTR^.NEXT <> NIL DO TOKENPTR := TOKENPTR^.NEXT;
767     NEW(TOKENPTR^.NEXT);
768     TOKENPTR := TOKENPTR^.NEXT;
769     GETTOKEN;
770     TOKENPTR^.TOKENVAL := TOKEN;
771     TOKENPTR^.NEXT := NIL
772     END;
773
774
775
776     PROCEDURE ERROR (ERRORCODE, DATA : INTEGER);
777
778     { Outputs an appropriate error message, and either obtains a replacement
779       TOKEN, or calls FAILURE. }
780
781     VAR TOKENLENGTH,
782         M           : INTEGER;
783         TOKENPTR    : PTOKENLIST;
784
785     BEGIN
786     FOR M := 1 TO (13 + 7*ORD(DIAGNOSTICS) + VALIDLENGTH) DO WRITE(' ');
787     TOKENLENGTH := N - VALIDLENGTH - 1;
788     FOR M := 1 TO TOKENLENGTH DO WRITE('^');
789     WRITELN;
790     WRITELN('**** ERROR',ERRORCODE : 2);
791     WRITEMESSAGE(ERRORCODE, DATA, ' ');
792     WRITELN;
793     CASE INPUTMODE OF
794     STOREDGENSAL : FAILURE(40, 0, ' ');
795     INSERTTEXT   : PROGERROR(2); {Error in SPSFILE expression}
796     TERMINAL     : BEGIN
797                     WRITELN('Remainder of input line ignored');
798                     WRITELN;
799                     FOR M := (VALIDLENGTH + 1) TO MAXLENGTH DO CURRENTLINE^.LINE[M] := ' ';
800                     N := MAXLENGTH;

```

```

801         GETTOKEN;
802         IF TOKENLIST <> NIL THEN
803             BEGIN
804                 {Need to put this token at the bottom of the list to
805                 replace the erroneous one. }
806                 TOKENPTR := TOKENLIST;
807                 WHILE TOKENPTR^.NEXT <> NIL DO TOKENPTR := TOKENPTR^.NEXT;
808                 TOKENPTR^.TOKENVAL := TOKEN
809             END
810         END
811     END
812 END;
813
814
815
816 FUNCTION MAGNITUDE(BOND : BONDORDER) : INTEGER;
817
818 { Returns an integer between 1 and 3 for the magnitude of BOND.
819   Called from PROCESSCT\HNUMBER
820   GETAVAILABLEPOSITIONS\MINPARENTBOND
821   GETAVAILABLEPOSITIONS\SUMFILIALS
822   PROCESSCT\GETPOSNS\GETSETPOSNS
823   PROCESSCT\GETPOSNS
824   ALTNVLIST\ELEMENT\SETCOMBARS\NEEDTOCHECK
825   ALTNVLIST\ELEMENT\SETCOMBARS\COMBINEDPOSITIONS
826   ALTNVLIST\ELEMENT\SETCOMBARS
827   ALTNVLIST\GETCHILDPOSITIONS
828   ALTNVLIST\PPOSNS\LMAGNOCHECKS
829   ASSIGNMENTSTMNT\POINTERLIST\ADDCOMBSUBS}
830
831
832 BEGIN
833 CASE BOND OF
834     NOTSPECIFIED,
835     ANY, SINGLE,
836     CHAIN, CHAISING,
837     RING, RINGSING      : MAGNITUDE := 1;
838     DOUBLE, CHAIDOUB,
839     RINGDOUB, AROMATIC,
840     RINGTAUT, CHAITAUT  : MAGNITUDE := 2;

```

```

841     TRIPLE,
842     CHAIRIP, RINGTRIP : MAGNITUDE := 3
843     END
844 END;
845
846
847
848 {-----}
849 PROCEDURE GETAVAILABLEPOSITIONS(PTRPS      : PTRPSTYPE;
850                                VAR POSNS  : INTEGSET;
851                                BOND MAG   : TBOND MAG);
852
853 { Returns in POSNS those positions of PTRPS^ which are substitutable, having a
854   sufficient number of spare valencies to accomodate a bond of magnitude BOND MAG.
855   Called by PROCESSCT\GETPOSNS
856           ALTNVLIST\ELEMENT\TRANSLATENOMEN\MODIFYGATEPOSITIONS
857           ALTNVLIST\ELEMENT\SETCOMBARS\COMBINEDPOSITIONS
858           ALTNVLIST\ELEMENT\PARAMETERLIST\FINDCONNECTIONS
859           ALTNVLIST\ALTNTVE\ADDPARALT
860           ALTNVLIST\PPOSNS
861           ASSIGNMENTSTMNT\POINTERLIST\ADDCOMBSUBS}
862
863 VAR ROWNO      : ATOMNUMBER;
864
865
866
867 FUNCTION MINBOND(OLDMAG,
868                NEWMAG : TBOND MAG) : TBOND MAG;
869
870 { Returns the smaller of the two values passed as parameter }
871
872 BEGIN
873 IF NEWMAG < OLDMAG
874     THEN MINBOND := NEWMAG
875     ELSE MINBOND := OLDMAG
876 END;
877
878
879
880 FUNCTION MINPARENTBOND(PARENTGATE : PPARENTLIST;

```

```

881             ROWNO      : ATOMNUMBER): TBOND MAG;
882
883 { Returns the magnitude of the smallest BOND to a parent if all items in the
884 list have ROWNO as the only element of CHILDPOSITIONS (i.e. there are no
885 alternatives). Otherwise, or if there is no parent list (PARENTGATE=NIL), returns 0.}
886
887 VAR VALID      : BOOLEAN;
888     COMBPOSNS  : PDOUBLIST;
889     MINPB      : TBOND MAG;
890
891 BEGIN
892     VALID := PARENTGATE <> NIL;    {initialisation}
893     MINPB := 3;    {initialise to large value}
894     WHILE VALID AND (PARENTGATE<>NIL) DO WITH PARENTGATE^ DO
895     BEGIN
896         IF CHILDPOSITIONS.COMBINED
897         THEN BEGIN
898             VALID := CHILDPOSITIONS.COMBMEMS <> NIL;
899             COMBPOSNS := CHILDPOSITIONS.COMBMEMS;
900             IF CONNBONDS.CONNECTIONS <> 2
901             THEN PROGERROR(3); {Combined childpositions with connections <> 2}
902             WHILE VALID AND (COMBPOSNS<>NIL) DO WITH COMBPOSNS^ DO
903             BEGIN
904                 WITH CONNBONDS DO
905                     IF (ROWNO=FIRST) AND (ROWNO=SECOND)
906                     THEN MINPB := MINBOND(MINPB, MAGNITUDE(BONDA) + MAGNITUDE(BONDB))
907                     ELSE IF ROWNO = FIRST
908                     THEN MINPB := MINBOND(MINPB, MAGNITUDE(BONDA))
909                     ELSE IF ROWNO = SECOND
910                     THEN MINPB := MINBOND(MINPB, MAGNITUDE(BONDB))
911                     ELSE VALID := FALSE;
912                 COMBPOSNS := NEXT
913             END
914         END
915     ELSE IF [ROWNO] = CHILDPOSITIONS.MEMBERS
916     THEN BEGIN
917         IF CONNBONDS.CONNECTIONS <> 1
918         THEN PROGERROR(4); {Uncombined childpositions with CONNECTIONS <> 1}
919         MINPB := MINBOND(MINPB, MAGNITUDE(CONNBONDS.BOND))
920     END

```

APPENDIX 3:

GENSAL INTERPRETER

```

921             ELSE VALID := FALSE;
922             PARENTGATE := NEXT
923             END;
924             IF VALID THEN MINPARENTBOND := MINPB
925                 ELSE MINPARENTBOND := 0
926             END;
927
928
929
930             FUNCTION SUMFILIALS(CONGENERS : CONGARRAY) : INTEGER;
931
932             { Returns the sum of the MAGNITUDES of FILIAL bonds }
933
934             VAR CNGNR : 1..MAXCONGENERS;
935                 SF : INTEGER;
936
937             BEGIN
938             SF := 0;
939             FOR CNGNR := 1 TO MAXCONGENERS DO WITH CONGENERS[CNGNR] DO
940                 IF RELATIONSHIP = FILIAL
941                     THEN SF := SF + MAGNITUDE(BOND);
942             SUMFILIALS := SF
943             END;
944
945
946
947             BEGIN {Body of GETAVAILABLEPOSITIONS}
948             CASE PTRPS^.PSVARIETY OF
949             DUMMY,
950             UNKNOWN,
951             OTHER : POSNS := [1..MAXCT];
952             GENERIC : WITH PTRPS^, PARAMLIST[ATOMCOUNT] DO
953                 IF TOPRANGE = NOTSET
954                     THEN IF SUBRANGES = NIL
955                         THEN POSNS := []
956                             ELSE POSNS := [1..SUBRANGES^.SECOND]
957                     ELSE POSNS := [1..MAXCT];
958             SPECIFIC : BEGIN
959                 POSNS := [];
960                 FOR ROWNO := 1 TO MAXCT DO IF PTRPS^.CT[ROWNO] <> NIL

```

```

961         THEN IF PTRPS^.CT[ROWNO]^ .ATOMICROW
962             THEN IF (PTRPS^.CT[ROWNO]^ .HYDROGENS
963                 - MINPARENTBOND(PTRPS^.PARENTGATE, ROWNO)
964                 - SUMFILIALS(PTRPS^.CT[ROWNO]^ .CONGENERSS))
965                 >= BOND MAG
966                 THEN POSNS := POSNS + [ROWNO]
967             END
968         END
969     END; { of GETAVAILABLEPOSITIONS
970     -----}
971
972
973
974     PROCEDURE LISTPOSNS (VAR LISTPTR : PDOUBLIST;
975         POSNSETA,
976         POSNSETB,
977         COMBPOSNS : INTEGSET);
978
979     { Returns a linked list of pairs of positions, being all the possible
980     combinations of the positions in POSNSETA and POSNSETB. The values in
981     any one item may only be identical if that value is in COMBPOSNS
982     Called by PROCESSCT\GETPOSNS
983     ALTNVLIST\GETCHILDPOSITIONS
984     ALTNVLIST\MODIFYCHILDPOSITIONS\GETCOMBPOSNS
985     ALTNVLIST\ELEMENT\SETCOMBARS\COMBINEDPOSITIONS
986     ALTNVLIST\ELEMENT\EXTRALAYER
987     ALTNVLIST\ELEMENT\GETLIMITPOSITIONS
988     ASSIGNMENTSTMNT\POINTERLIST\ADDCOMBSUBS
989     ALTNVLIST\ELEMENT\PARAMETERLIST\FINDCONNECTIONS}
990
991     VAR POSNA,
992         POSNB : INTEGER;
993         NEWITEM : PDOUBLIST;
994
995     BEGIN
996     FOR POSNA := 0 TO MAXCT DO
997         IF POSNA IN POSNSETA
998             THEN FOR POSNB := 0 TO MAXCT DO
999                 IF (POSNB IN POSNSETB) AND ((POSNB<>POSNA) OR (POSNB IN COMBPOSNS))
1000                 THEN BEGIN

```

```

1001                 NEW(NEWITEM);
1002                 ECTRSIZE := ECTRSIZE + 6;
1003                 NEWITEM^.FIRST := POSNA;
1004                 NEWITEM^.SECOND := POSNB;
1005                 NEWITEM^.NEXT := LISTPTR;
1006                 LISTPTR := NEWITEM
1007                 END
1008     END;
1009
1010
1011
1012     PROCEDURE ADDTOLIST(NEWPS : PTRPSTYPE);
1013
1014     { Adds a PS to the bottom of the list of PSs.
1015       Called by COPYPS
1016         ALTNVLIST\ELEMENT\SETCOMBARS
1017         ALTNVLIST\ELEMENT\SUBSTASVALUE}
1018
1019     VAR NEWIRITEM : PIRLIST;
1020
1021     BEGIN
1022     NEW(NEWIRITEM);
1023     ECTRSIZE := ECTRSIZE + 4;
1024     NEWIRITEM^.PARSTRUCT := NEWPS;
1025     NEWIRITEM^.NEXT := NIL;
1026     IRLISTBOT^.NEXT := NEWIRITEM;
1027     IRLISTBOT := NEWIRITEM
1028     END;
1029
1030
1031
1032     FUNCTION COPYPS(OLDPS : PTRPSTYPE) : PTRPSTYPE;
1033
1034     { Copies a PS.
1035       Called by COPYCOMBAR
1036         ALTNVLIST\ELEMENT\SUBSTASVALUE
1037         ALTNVLIST\ELEMENT\SETCOMBARS}
1038
1039     VAR NEWPS : PTRPSTYPE;
1040

```

```

1041 BEGIN
1042 CASE OLDPS^.PSVARIETY OF
1043   DUMMY   : BEGIN
1044     NEW(NEWPS, DUMMY);
1045     ECTRSIZE := ECTRSIZE + 8;
1046     NEWPS^.SUBSTNAME := OLDPS^.SUBSTNAME
1047   END;
1048   UNKNOWN : BEGIN
1049     NEW(NEWPS, UNKNOWN);
1050     ECTRSIZE := ECTRSIZE + 6
1051   END;
1052   OTHER   : BEGIN
1053     NEW(NEWPS, OTHER);
1054     ECTRSIZE := ECTRSIZE + 22;
1055     NEWPS^.TERM := OLDPS^.TERM
1056   END;
1057   SPECIFIC : BEGIN
1058     NEW(NEWPS, SPECIFIC);
1059     ECTRSIZE := ECTRSIZE + 70;
1060     NEWPS^.CT := OLDPS^.CT
1061   END;
1062   GENERIC  : BEGIN
1063     NEW(NEWPS, GENERIC);
1064     ECTRSIZE := ECTRSIZE + 50;
1065     NEWPS^.PARAMLIST := OLDPS^.PARAMLIST
1066   END
1067 END;
1068 NEWPS^.VISITED := FALSE;
1069 NEWPS^.PARENTGATE := NIL;
1070 NEWPS^.CHILDGATE := NIL;
1071 NEWPS^.PSVARIETY := OLDPS^.PSVARIETY;
1072 ADDTOLIST(NEWPS);
1073 COPYPS := NEWPS
1074 END;
1075
1076
1077
1078 PROCEDURE COPYCOMBAR(VAR NEWCOMBAR : PCOMBINLIST;
1079                     OLDCOMBAR,
1080                     LASTCOMBLAYER : PCOMBINLIST;

```

```

1081                 LASTPPOSNS      : PTGROUPMEMS;
1082                 PRNTPS           : PTRPSTYPE;
1083                 FIRSTBAR,
1084                 OMITPG,
1085                 COPYPSS           : BOOLEAN);
1086 FORWARD;
1087
1088
1089
1090 PROCEDURE COPYALTBAR(VAR NEWALTBAR : PALTERNLIST;
1091                    OLDCOMBLIST,
1092                    LASTCOMBLAYER : PCOMBINLIST;
1093                    LASTPPOSNS    : PTGROUPMEMS;
1094                    PRNTPS        : PTRPSTYPE;
1095                    FIRSTBAR,
1096                    OMITPG,
1097                    COPYPSS      : BOOLEAN);
1098
1099 { Creates a new alternative bar item, and calls COPYCOMBAR to copy the
1100   combination bar items in OLDCOMBLIST into its COMBINATION field.
1101   Called by COPYCOMBAR
1102     ALTNVLIST\ELEMENT\SUBSTASVALUE}
1103 VAR NEWAB : PALTERNLIST;
1104
1105 BEGIN
1106   NEW(NEWAB);
1107   ECTRSIZE := ECTRSIZE + 4;
1108   NEWAB^.COMBINATION := NIL;
1109   WHILE OLDCOMBLIST <> NIL DO
1110     BEGIN
1111       COPYCOMBAR(NEWAB^.COMBINATION, OLDCOMBLIST, LASTCOMBLAYER, LASTPPOSNS, PRNTPS, FIRSTBAR, OMITPG, COPYP
1112     S);
1113       OLDCOMBLIST := OLDCOMBLIST^.NEXT
1114     END;
1115   NEWAB^.NEXT := NEWALTBAR;
1116   NEWALTBAR := NEWAB
1117 END;
1118
1119

```

APPENDIX 3:

GENSAL INTERPRETER

```
1120 PROCEDURE SETCONNBONDS(VAR CONNBONDS : TCONNBONDS;
1121                          CONNECTIVITY : TCONNS);
1122
1123 { Returns a TCONNBONDS record with CONNECTIONS set to CONNECTIVITY and all
1124   bond types to NOTSPECIFIED.
1125   Called by SUBSTGROUP\CHECKCOMPATIBILITY
1126           ALTNVLIST\ELEMENT\VALIDSUBST
1127           ALTNVLIST\ELEMENT
1128           ASSIGNMENTSTMNT\POINTERLIST\ADDDEFNTABLE
1129   DECLARESUBST}
1130
1131 BEGIN
1132 WITH CONNBONDS DO
1133   BEGIN
1134     CONNECTIONS := CONNECTIVITY;
1135     CASE CONNECTIONS OF
1136     NOTSET,
1137     0       : ;
1138     1       : BOND := NOTSPECIFIED;
1139     2       : BEGIN
1140               BONDA := NOTSPECIFIED;
1141               BONDB := NOTSPECIFIED
1142             END
1143     END
1144   END
1145 END;
1146
1147
1148
1149 PROCEDURE UPDATEPPSCONNS(PARPSLIST : PPSLIST);
1150
1151 { Copies the CONNBONDS field of the first item in PARPSLIST right down the
1152   list.
1153   Called by SUBSTGROUP\CHECKCOMPATIBILITY
1154           ALTNVLIST\ELEMENT\VALIDSUBST
1155   DECLARESUBST}
1156
1157 VAR NEWCONNBONDS : TCONNBONDS;
1158
1159 BEGIN
```

```

1160     NEWCONN BONDS := PARPSLIST^.CONN BONDS;
1161     REPEAT
1162         PARPSLIST^.CONN BONDS := NEWCONN BONDS;
1163         PARPSLIST := PARPSLIST^.NEXT
1164     UNTIL PARPSLIST = NIL
1165     END;
1166
1167
1168
1169     PROCEDURE DECLARESUBST(SUBST          : SUBSTITUENT;
1170                          PSADDRESS,
1171                          SAVPS          : PTRPSTYPE;
1172                          CONN BONDS    : TCONN BONDS;
1173                          PRNTPOSNS     : PTGROUPMEMS);
1174
1175     { Adds declaration of SUBST to RDECLARATIONTABLE, referencing PSADDRESS.
1176       If this is the first declaration of SUBST to give a value for connectivity,
1177       then the connectivity is copied into all the other declarations.
1178       Called by PROCESSCT
1179         COPYCOMBAR
1180         ALTNVLIST\ELEMENT\PARAMETERLIST\USERPARAMETER
1181         ALTNVLIST\ELEMENT\SUBSTASVALUE}
1182
1183     VAR PTR : PPSLIST;
1184
1185     BEGIN
1186     IF RDECLARATIONTABLE[SUBST] <> NIL
1187         THEN IF (RDECLARATIONTABLE[SUBST]^ .CONN BONDS.CONNECTIONS = NOTSET)
1188             AND (CONN BONDS.CONNECTIONS <> NOTSET)
1189             THEN BEGIN
1190                 SETCONN BONDS(RDECLARATIONTABLE[SUBST]^ .CONN BONDS, CONN BONDS.CONNECTIONS);
1191                 UPDATEPPSCONNS(RDECLARATIONTABLE[SUBST])
1192             END;
1193     NEW(PTR);
1194     PTR^.PARSTRUCT := PSADDRESS;
1195     PTR^.COMBINS   := NIL;
1196     PTR^.FURTHERSUB := SAVPS;
1197     PTR^.CONN BONDS := CONN BONDS;
1198     PTR^.PRNTPOSNS := PRNTPOSNS;
1199     PTR^.COPYCHILDPS := FALSE;

```

```

1200 PTR^.NEXT := RDECLARATIONTABLE[SUBST];
1201 RDECLARATIONTABLE[SUBST] := PTR;
1202 DECLSUBS := DECLSUBS + [SUBST]
1203 END;
1204
1205
1206
1207 PROCEDURE COMPARELISTS(LOWERLIST,
1208                        UPPERLIST : PDOUBLIST);
1209
1210 { Calls FAILURE if the items in LOWERLIST are not identical with those
1211   in UPPERLIST.
1212   Called by CHECKINCLUDED}
1213
1214 VAR PTR      : PDOUBLIST;
1215     FOUND    : BOOLEAN;
1216
1217 BEGIN
1218   WHILE LOWERLIST <> NIL DO
1219     BEGIN
1220       PTR := UPPERLIST;
1221       WHILE PTR <> NIL DO
1222         BEGIN
1223           FOUND := (LOWERLIST^.FIRST = PTR^.FIRST) AND (LOWERLIST^.SECOND = PTR^.SECOND);
1224           IF FOUND THEN PTR := NIL
1225             ELSE PTR := PTR^.NEXT
1226         END;
1227       IF FOUND THEN LOWERLIST := LOWERLIST^.NEXT
1228         ELSE FAILURE(39, 0, ' ')
1229     END
1230   END;
1231
1232
1233
1234 PROCEDURE CHECKALLWITHIN(COBBMEMS : PDOUBLIST;
1235                          AVAILPOSNS : INTEGSET;
1236                          FAILCODE : INTEGER);
1237
1238 { Checks that all the items in COBBMEMS are within AVAILPOSNS
1239   Called by ELEMENT\TRANSLATENOMEN\MODIFYGATEPOSITIONS

```

```

1240             ALTNVLIST\POSNS\LMAGCHECKS
1241             CHECKINCLUDED}
1242
1243 BEGIN
1244 WHILE COMBMEMS <> NIL DO WITH COMBMEMS^ DO
1245     BEGIN
1246         IF [FIRST, SECOND] <= AVAILPOSNS
1247             THEN {OK}
1248             ELSE IF FIRST IN AVAILPOSNS
1249                 THEN FAILURE(FAILCODE, SECOND, '  ')
1250                 ELSE FAILURE(FAILCODE, FIRST, '  ');
1251     COMBMEMS := NEXT
1252     END
1253 END;
1254
1255
1256
1257 PROCEDURE CHECKINCLUDED(LOWERPOSNS,
1258                         UPPERPOSNS : PTGROUPMEMS);
1259
1260 { Checks that all the positions in LOWERPOSNS are also in UPPERPOSNS
1261   Called by COPYCOMBAR}
1262
1263 BEGIN
1264 IF LOWERPOSNS^.COMBINED
1265     THEN IF UPPERPOSNS^.COMBINED
1266         THEN COMPARELISTS(LOWERPOSNS^.COMBMEMS, UPPERPOSNS^.COMBMEMS)
1267         ELSE CHECKALLWITHIN(LOWERPOSNS^.COMBMEMS, UPPERPOSNS^.MEMBERS, 39)
1268     ELSE IF UPPERPOSNS^.COMBINED
1269         THEN PROGERROR(26) {Trying to uncombine a position set}
1270         ELSE IF LOWERPOSNS^.MEMBERS <= UPPERPOSNS^.MEMBERS
1271             THEN {OK}
1272             ELSE FAILURE(39, 0, '  ')
1273 END;
1274
1275
1276
1277
1278 PROCEDURE COPYCOMBAR; {Previous FORWARD declaration}
1279

```

```

1280 { Copies a combination bar item. If FIRSTBAR then the PARENTPOSITIONS field
1281 is altered to LASTPPOSNS. In other cases LASTPPOSNS is changed to
1282 OLDCOMBAR^.PARENTPOSITIONS after checking that the specified positions are
1283 available. For BOTTOMBARs a new parent gate is created on the existing
1284 Child PS (not done if OMITPG). If COPYPSS is TRUE, then Child PSs
1285 themselves are copied, otherwise the new Gate is simply made to point to the
1286 Child PS. For non-BOTTOMBARs, COPYALTBAR is called to copy the ALTERNATIVES.
1287 Called by COPYALTBAR
1288 COPYCOMBAR (recursively)
1289 ENTERCOMBIN
1290 ALTNVLIST\ADDFURTHERSUBTN}
1291
1292 VAR NEWCB,
1293     SUBCB      : PCOMBINLIST;
1294     OLDALTBAR : PALTERNLIST;
1295     NEWPG     : PPARENTLIST;
1296
1297
1298
1299 BEGIN
1300 IF OLDCOMBAR^.BOTTOMBAR
1301     THEN NEW(NEWCB, TRUE)
1302     ELSE NEW(NEWCB, FALSE);
1303 ECTRSIZE := ECTRSIZE + 11 + ORD(OLDCOMBAR^.BOTTOMBAR) * 13;
1304 IF FIRSTBAR
1305     THEN NEWCB^.PARENTPOSITIONS := LASTPPOSNS
1306     ELSE IF OLDCOMBAR^.PARENTPOSITIONS = NIL
1307         THEN NEWCB^.PARENTPOSITIONS := NIL
1308         ELSE BEGIN
1309             NEWCB^.PARENTPOSITIONS := OLDCOMBAR^.PARENTPOSITIONS;
1310             CHECKINCLUDED(NEWCB^.PARENTPOSITIONS, LASTPPOSNS);
1311             LASTPPOSNS := NEWCB^.PARENTPOSITIONS
1312         END;
1313 NEWCB^.FREQUENCY := OLDCOMBAR^.FREQUENCY;
1314 NEWCB^.NEXT := NEWCOMBAR;
1315 NEWCB^.BOTTOMBAR := OLDCOMBAR^.BOTTOMBAR;
1316 IF NEWCB^.BOTTOMBAR
1317     THEN BEGIN
1318         IF COPYPSS
1319             THEN NEWCB^.CHILDPs := COPYPS(OLDCOMBAR^.CHILDPs)

```

```

1320 ELSE NEWCB^.CHILDPS := OLDCOMBAR^.CHILDPS;
1321 IF NEWCB^.CHILDPS^.PSVARIETY = DUMMY
1322 THEN BEGIN
1323     DECLARESUBST(NEWCB^.CHILDPS^.SUBSTNAME,
1324                 PRNTPS,
1325                 NEWCB^.CHILDPS,
1326                 OLDCOMBAR^.CONNBONDS,
1327                 LASTPPOSNS);
1328     RDECLARATIONTABLE[NEWCB^.CHILDPS^.SUBSTNAME]^COMBINS := LASTCOMBLAYER
1329     END;
1330 NEWCB^.CHILDPOSITIONS := OLDCOMBAR^.CHILDPOSITIONS;
1331 NEWCB^.CONNBONDS := OLDCOMBAR^.CONNBONDS;
1332 IF NOT OMITPG
1333 THEN BEGIN
1334     NEW(NEWPG);
1335     ECTRSIZE := ECTRSIZE + 26;
1336     WITH NEWPG^ DO
1337     BEGIN
1338         CHILDPOSITIONS := NEWCB^.CHILDPOSITIONS;
1339         PARENTPOSITIONS := LASTPPOSNS^;
1340         PARENTPS := PRNTPS;
1341         CONNBONDS := NEWCB^.CONNBONDS;
1342         NEXT := NEWCB^.CHILDPS^.PARENTGATE
1343     END;
1344     NEWCB^.CHILDPS^.PARENTGATE := NEWPG
1345     END;
1346 IF COPYPSS
1347 THEN BEGIN
1348     SUBCB := OLDCOMBAR^.CHILDPS^.CHILDGATE;
1349     WHILE SUBCB <> NIL DO
1350     BEGIN
1351         COPYCOMBAR(NEWCB^.CHILDPS^.CHILDGATE, SUBCB, NIL, NIL, NEWCB^.CHILDP
1352     TRUE);
1353         SUBCB := SUBCB^.NEXT
1354     END
1355     END
1356 ELSE BEGIN
1357     OLDALTBAR := OLDCOMBAR^.ALTERNATIVES;
1358     NEWCB^.ALTERNATIVES := NIL;

```

APPENDIX 3:

GENSAL INTERPRETER

FALSE, OMITPG

```

1359         WHILE OLDALTBAR <> NIL DO
1360             BEGIN
1361                 COPYALTBAR(NEWCB^.ALTERNATIVES, OLDALTBAR^.COMBINATION, NEWCB,
1362                     LASTPPOSNS, PRNTPS, FALSE, OMITPG, COPYPSS);
1363                 OLDALTBAR := OLDALTBAR^.NEXT
1364             END
1365         END;
1366     NEWCOMBAR := NEWCB
1367 END;

1368
1369
1370
1371     PROCEDURE ENTERCOMBIN(SUBST          : SUBSTITUENT;
1372                          VAR GATEENTRY  : PCOMBINLIST);
1373
1374     { If SUBST has been defined, copies the existing definition combination bar
1375       into GATEENTRY, otherwise creates a new non-BOTTOMBAR combination bar item.
1376       Called by PROCESSCT
1377         ALTNVLIST\ELEMENT\PARAMETERLIST\USERPARAMETER}
1378
1379     BEGIN
1380     WITH RDECLARATIONTABLE[SUBST]^ DO
1381     IF RDEFINITIONTABLE[SUBST] = NIL
1382     THEN BEGIN
1383         NEW(COMBINS, TRUE);
1384         ECTRSIZE := ECTRSIZE + 11;
1385         COMBINS^.BOTTOMBAR := FALSE;
1386         COMBINS^.PARENTPOSITIONS := PRNTPOSNS;
1387         COMBINS^.FREQUENCY.TOPRANGE := NOTSET;
1388         COMBINS^.FREQUENCY.SUBRANGES := ESSENTFREQ;
1389         COMBINS^.ALTERNATIVES := NIL;
1390         COMBINS^.NEXT := GATEENTRY;
1391         GATEENTRY := COMBINS
1392     END
1393     ELSE BEGIN
1394         COPYCOMBAR(GATEENTRY,
1395                 RDEFINITIONTABLE[SUBST],
1396                 NIL,
1397                 PRNTPOSNS,
1398                 PARSTRUCT,

```

```

1399             TRUE,
1400             FALSE,
1401             FALSE);
1402     COMBINS := GATEENTRY
1403     END
1404 END;
1405
1406
1407
1408
1409 {-----}
1410             PROCEDURE PROCESSCT
1411 -----}
1412
1413
1414
1415 PROCEDURE PROCESSCT (VAR CT      : CTYPE;
1416                    INTERACTIVE : BOOLEAN;
1417                    PSADDRESS   : PTRPSTYPE);
1418
1419 { Carries out the reformatting of FELDCT/FELDBD, putting the result into CT.
1420   Called by READSD
1421   ALTNVLIST\ELEMENT\TRANSLATENOMEN}
1422
1423 TYPE TNODENATURE = (ATOMIC, APICLABEL, VARPOSNLABEL, EXPHYDROGEN, SUBSTNODE);
1424
1425 VAR APICCOUNT,           { Number of apical labels present }
1426     ROWNO      : ATOMNUMBER; { Row counter for CT and FELDCT }
1427     M          : 1..2;      { Counter for characters of atomic symbol }
1428     REJECTED   : BOOLEAN;   { Set to TRUE if error is found }
1429     CONNBONDS : TCONNBONDS;
1430     PRNTPOSNS : PTGROUPMEMS;
1431
1432
1433
1434 FUNCTION NODENATURE(NODE : ATOMNUMBER) : TNODENATURE;
1435
1436 { Returns the nature of the NODE in FELDCT.
1437   Called by INDEPENDENT
1438   READCONGNERS

```

```
1439             Body of PROCESSCT}
1440
1441 BEGIN
1442 WITH FELDCT[NODE] DO
1443 IF CHEM = '*'
1444     THEN NODENATURE := APICLABEL
1445     ELSE IF CHEM = '#
1446         THEN NODENATURE := VARPOSNLABEL
1447         ELSE IF CHEM = 'H
1448             THEN NODENATURE := EXPHYDROGEN
1449             ELSE IF CHEM[2] IN ['0'..'9']
1450                 THEN NODENATURE := SUBSTNODE
1451                 ELSE NODENATURE := ATOMIC
1452 END;
1453
1454
1455
1456 FUNCTION BONDVAL (NODEA,NODEB : ATOMNUMBER) : BONDORDER;
1457 { Finds the order of the bond between NODEA and NODEB }
1458
1459 VAR M      : 0..MAXCT;
1460     BNDVAL : BONDORDER;
1461
1462 BEGIN
1463 BNDVAL := NOTSPECIFIED;
1464 M := 0;
1465 WHILE (M <= NUMOFBONDS) AND (BNDVAL=NOTSPECIFIED) DO
1466 BEGIN
1467     M := M+1;
1468     WITH FELDBD[M] DO
1469         IF ((NODEA=NODE1) AND (NODEB=NODE2))
1470         OR ((NODEA=NODE2) AND (NODEB=NODE1))
1471         THEN REPEAT BNDVAL := SUCC(BNDVAL)
1472             UNTIL ORD(BNDVAL) = BOND
1473     END;
1474 BNDVAL := BNDVAL
1475 END;
1476
1477
1478
```

```

1479
1480 FUNCTION SUBSTNAME(CHEM : STRING4): SUBSTITUENT;
1481
1482 { Converts the name of a substituent to integer format from characters.
1483   Called by Body of PROCESSCT}
1484
1485 VAR SUBST : SUBSTITUENT;
1486
1487 BEGIN
1488 IF CHEM[3] = ' '
1489 THEN SUBST := ORD(CHEM[2]) - ORD('0')
1490 ELSE SUBST := (ORD(CHEM[3]) - ORD('0')) + 10*(ORD(CHEM[2]) - ORD('0')) ;
1491 SUBSTNAME := SUBST
1492 END;
1493
1494
1495
1496 FUNCTION INDEPENDENT(NODENO : ATOMNUMBER): BOOLEAN;
1497
1498 { Returns TRUE if any of the congeners of NODENO are EXPHYDROGEN, APICLABEL
1499   or SUBSTNODE.
1500   Called by Body of PROCESSCT}
1501
1502 VAR CNGNR : 1..MAXCONGENERS;
1503
1504 BEGIN
1505 INDEPENDENT := FALSE;
1506 FOR CNGNR := 1 TO MAXCONGENERS DO IF FELDCT[NODENO].AR[CNGNR] <> 0
1507 THEN IF NODENATURE(FELDCT[NODENO].AR[CNGNR]) IN [EXPHYDROGEN, APICLABEL, SUBSTNODE]
1508 THEN INDEPENDENT := TRUE
1509 END;
1510
1511
1512 PROCEDURE REJECT(ERRORCODE : INTEGER;
1513                 NODE       : ATOMNUMBER);
1514
1515 { Outputs an error message.
1516   Called by READCONGENERS
1517     HNUMBER
1518     NUMOFCONNS

```

```

1519             CHECKEARLIERDEFN
1520             Body of PROCESSCT}
1521
1522 BEGIN
1523 WRITE('SD ERROR: ');
1524 WRITEMESSAGE(ERRORCODE, NODE, ' ');
1525 REJECTED := TRUE
1526 END;
1527
1528
1529
1530 PROCEDURE READCONGENERS (VAR CONGENERS : CONGARRAY;
1531                          VAR HYDROGENS : NUMCONGENERS;
1532                          ATOMICROW    : BOOLEAN;
1533                          ROWNO        : ATOMNUMBER);
1534
1535 { Sets the values in CONGENERS, and the number
1536   of explicit attached HYDROGENS, for a single connection table ROW.
1537   Called by Body of PROCESSCT}
1538
1539 VAR FELDCONG : ARRAY[1..MAXCONGENERS] OF ATOMNUMBER;
1540     CNGNR    : 1..MAXCONGENERS;
1541
1542 BEGIN
1543 HYDROGENS := 0;
1544 FELDCONG := FELDCT[ROWNO].AR;
1545 FOR CNGNR := 1 TO MAXCONGENERS DO
1546   BEGIN
1547     CONGENERS[CNGNR].RELATIONSHIP := NONE;
1548     IF FELDCONG[CNGNR] <> 0 THEN WITH CONGENERS[CNGNR-HYDROGENS] DO
1549       BEGIN
1550         BOND := BONDVAL(ROWNO, FELDCONG[CNGNR]);
1551         CASE NODENATURE(FELDCONG[CNGNR]) OF
1552           ATOMIC      : BEGIN
1553             RELATIONSHIP := FRATERNAL;
1554             ROWNUM := FELDCONG[CNGNR]
1555           END;
1556           EXPHYDROGEN : HYDROGENS := HYDROGENS + 1;
1557           APICLABEL   : RELATIONSHIP := PARENTAL;
1558           VARPOSNLABEL : IF ATOMICROW

```

```

1559         THEN REJECT(56, ROWNO)
1560         ELSE BEGIN
1561             RELATIONSHIP := FRATERNAL;
1562             ROWNUM := NOTFIXED
1563         END;
1564     SUBSTNODE : IF INDEPENDENT(FELDCONG[CNGNR])
1565     THEN BEGIN
1566         RELATIONSHIP := FRATERNAL;
1567         ROWNUM := FELDCONG[CNGNR]
1568     END
1569     ELSE RELATIONSHIP := FILIAL
1570     END
1571     END
1572     END
1573 END;
1574
1575
1576
1577 PROCEDURE HNUMBER (NODE : ATOMNUMBER);
1578
1579 { Sets a value for HYDROGENS at atom NODE, checking valencies in VELENCYFILE.
1580   Called by Body of PROCESSCT}
1581
1582 CONST MAXSTATES = 5;
1583
1584 TYPE TELEMVALS = RECORD
1585     ELEMENT : STRING2;
1586     VALENCIES : PACKED ARRAY[1..MAXSTATES] OF NUMCONGENERS
1587 END;
1588
1589 VAR BONDCOUNT : 0..18;           { Sum of bond orders }
1590     STATE : INTEGER;
1591     ARCOUNT,           { Number of aromatic bonds }
1592     TAUTCOUNT,         { Number of tautomeric bonds }
1593     CNGNR : 0..MAXCONGENERS; { Congner counter }
1594     EXTERNBONDS,      { Sum of MAGNITUDEs of external bonds }
1595     SPAREVALS : INTEGER;      { Valency of common atom }
1596     ELEMVAL : TELEMVALS;      { Element valency record }
1597     VALENCYFILE : FILE OF TELEMVALS;
1598

```

```

1599 BEGIN
1600 BONDCount := 0;
1601 ARCount := 0;
1602 TAUTCount := 0;
1603 EXTERNBONDS := 0;
1604 FOR CNGNR := 1 TO MAXCONGENERES DO WITH CT[NODE]^CONGENERES[CNGNR] DO
1605 BEGIN
1606 IF RELATIONSHIP <> NONE THEN
1607 CASE BOND OF
1608 NOTSPECIFIED, RINGSING, CHAISING,
1609 RING, CHAIN, SINGLE, ANY : BONDCount := BONDCount +1;
1610 RINGDOUB, CHAIDOUB, DOUBLE : BONDCount := BONDCount +2;
1611 RINGTRIP, CHAITRIP, TRIPLE : BONDCount := BONDCount +3;
1612 AROMATIC : ARCount := ARCount +1;
1613 RINGTAUT, CHAITAUT : TAUTCount := TAUTCount +1
1614 END;
1615 IF RELATIONSHIP IN [FILIAL, PARENTAL]
1616 THEN EXTERNBONDS := EXTERNBONDS + MAGNITUDE(BOND)
1617 END;
1618 CASE ARCount OF
1619 0 : ;
1620 2,3 : BONDCount := BONDCount + ARCount +1;
1621 1,4,5,6 : REJECT(60, NODE)
1622 END;
1623 CASE TAUTCount OF
1624 0,1 : BONDCount := BONDCount+TAUTCount;
1625 2,3 : BONDCount := BONDCount + TAUTCount +1;
1626 4,5,6 : REJECT(61, NODE)
1627 END;
1628 RESET(VALENCYFILE, 'LI2GEN>VALENCYFILE');
1629 ELEMVAL.ELEMENT := ' ';
1630 WHILE(ELEMVAL.ELEMENT <> CT[NODE]^ATOM) AND NOT EOF(VALENCYFILE) DO
1631 READ(VALENCYFILE, ELEMVAL);
1632 IF EOF(VALENCYFILE)
1633 THEN {atom not in file}
1634 ELSE BEGIN
1635 STATE := 1;
1636 WHILE STATE <= MAXSTATES DO
1637 BEGIN
1638 SPAREVALS := ELEMVAL.VALENCIES[STATE] + CT[NODE]^CHARGE - BONDCount;

```

```
1639         IF SPAREVALS < 0
1640             THEN IF STATE < MAXSTATES
1641                 THEN STATE := STATE + 1
1642             ELSE BEGIN
1643                 REJECT(55, NODE);
1644                 STATE := MAXSTATES+1
1645             END
1646         ELSE IF SPAREVALS > 6
1647             THEN PROGERROR(5) {Excessively large valency}
1648         ELSE BEGIN
1649             CT[NODE]^.HYDROGENS := SPAREVALS + EXTERNBONDS;
1650             STATE := MAXSTATES + 1
1651         END
1652     END
1653 END;
1654 RESET(VALENCYFILE, '@TTY')
1655 END;
1656
1657
1658
1659 FUNCTION NUMOFCONNS(CONGENERES : CONGARRAY;
1660                    TOTCONNS   : NUMCONGENERES;
1661                    NODE       : ATOMNUMBER) : NUMCONGENERES;
1662
1663 { Returns the number of connections specified in CONGENERES, plus the entry
1664   value of TOTCONNS (which corresponds to the number of HYDROGENS.
1665   Called by CHECKEARLIERDEFN
1666   Body of PROCESSCT}
1667
1668 VAR CNGNR : NUMCONGENERES;
1669
1670 BEGIN
1671 FOR CNGNR := 1 TO MAXCONGENERES DO
1672     IF CONGENERES[CNGNR].RELATIONSHIP <> NONE
1673     THEN TOTCONNS := TOTCONNS + 1;
1674 IF TOTCONNS > 2 THEN REJECT(54, NODE);
1675 NUMOFCONNS := TOTCONNS
1676 END;
```

```

1679
1680 PROCEDURE CHECKEARLIERDEFN(SUBST      : SUBSTITUENT;
1681                          CONNS      : NUMCONGENERS;
1682                          ROWSREAD   : ATOMNUMBER);
1683
1684 { Examines the rows of the connection table up as far as ROWSREAD, and if a
1685   non-ATOMICROW is found with NAME=SUBST then compares the value of CONNS
1686   with the number of connections of this row.
1687   Called by Body of PROCESSCT}
1688
1689 VAR NODENO : ATOMNUMBER;
1690
1691 BEGIN
1692   NODENO := 1;
1693   WHILE NODENO < ROWSREAD DO
1694     IF CT[NODENO] = NIL
1695       THEN NODENO := NODENO + 1
1696       ELSE WITH CT[NODENO]^ DO
1697         IF ATOMICROW
1698           THEN NODENO := NODENO + 1
1699           ELSE IF NAME = SUBST
1700             THEN BEGIN
1701                   IF CONNS = NUMOFCONNS(CONGENERS, HYDROGENS, NODENO)
1702                     THEN {matches OK}
1703                     ELSE REJECT(53, ROWSREAD);
1704                   NODENO := ROWSREAD
1705                   END
1706                   ELSE NODENO := NODENO + 1
1707 END;
1708
1709
1710 {.....}
1711 PROCEDURE GETPOSNS(CONGENERS      : CONGARRAY;
1712                  VAR CONNBONDS   : TCONNBONDS;
1713                  VAR PRNTPOSNS   : PTGROUPMEMS);
1714
1715 { Sets CONNBONDS, PRNTPOSNS for a substituent, by examining CONGENERS.
1716   Called by Body of PROCESSCT}
1717
1718 VAR POSNS1,

```

```

1719         POSNS2,
1720         COMBPOSNS : INTEGSET;
1721         MAGSUM    : INTEGER;
1722         REVERSIBLE : BOOLEAN;
1723
1724
1725
1726     PROCEDURE GETSETPOSNS(VAR SETPOSNS : INTEGSET;
1727                          POSITION      : ATOMNUMBER;
1728                          BOND        : BONDORDER);
1729
1730     BEGIN
1731     IF POSITION = NOTFIXED
1732     THEN GETAVAILABLEPOSITIONS(PSADDRESS, SETPOSNS, MAGNITUDE(BOND))
1733     ELSE SETPOSNS := [POSITION]
1734     END;
1735
1736
1737
1738     BEGIN {Body of GETPOSNS}
1739     NEW(PRNTPOSNS);
1740     ECTRSIZE := ECTRSIZE + 9;
1741     IF CONGENERES[1].RELATIONSHIP = NONE
1742     THEN BEGIN
1743         {substituent is unconnected}
1744         PRNTPOSNS^.COMBINED := FALSE;
1745         PRNTPOSNS^.MEMBERS := [];
1746         CONNBONDS.CONNECTIONS := 0
1747     END
1748     ELSE IF CONGENERES[2].RELATIONSHIP = NONE
1749     THEN BEGIN
1750         {substituent is singly connected}
1751         PRNTPOSNS^.COMBINED := FALSE;
1752         CONNBONDS.CONNECTIONS := 1;
1753         CONNBONDS.BOND := CONGENERES[1].BOND;
1754         CASE CONGENERES[1].RELATIONSHIP OF
1755             FRATERNAL : WITH CONGENERES[1] DO
1756                 GETSETPOSNS(PRNTPOSNS^.MEMBERS, ROWNUM, BOND);
1757             FILIAL    : PROGERROR(6); {substituent node with filial congener }
1758             PARENTAL  : PRNTPOSNS^.MEMBERS := [0]

```

```

1759         END
1760     END
1761 ELSE BEGIN
1762     {substituent is doubly connected}
1763     PRNTPOSNS^.COMBINED := TRUE;
1764     PRNTPOSNS^.COMBMEMS := NIL;
1765     WITH CONNBONDS DO
1766     BEGIN
1767         CONNECTIONS := 2;
1768         BONDA := CONGENERES[1].BOND;
1769         BONDB := CONGENERES[2].BOND;
1770         REVERSIBLE := (BONDMATCHARRAY[BONDA,BONDB] IN [ANY, CHAIN, RING]) OR (BONDA=BONDB)
1771     END;
1772     CASE CONGENERES[1].RELATIONSHIP OF
1773     FRATERNAL :
1774     CASE CONGENERES[2].RELATIONSHIP OF
1775     FRATERNAL : BEGIN
1776         WITH CONGENERES[1] DO GETSETPOSNS(POSNS1, ROWNUM, BOND);
1777         WITH CONGENERES[2] DO GETSETPOSNS(POSNS2, ROWNUM, BOND);
1778         WITH CONNBONDS DO
1779             MAGSUM := MAGNITUDE(BONDA) + MAGNITUDE(BONDB);
1780             IF (POSNS1 * POSNS2 = []) OR (MAGSUM > 3)
1781             THEN COMBPOSNS := []
1782             ELSE BEGIN
1783                 GETAVAILABLEPOSITIONS(PSADDRESS, COMBPOSNS, MAGSUM);
1784                 COMBPOSNS := COMBPOSNS * POSNS1 * POSNS2
1785             END;
1786             LISTPOSNS(PRNTPOSNS^.COMBMEMS, POSNS1, POSNS2, COMBPOSNS);
1787             IF REVERSIBLE
1788             THEN LISTPOSNS(PRNTPOSNS^.COMBMEMS, POSNS2, POSNS1, COMBPOSNS)
1789             END;
1790     FILIAL      : PROGERROR(7); {substituent node with filial congener }
1791     PARENTAL   : BEGIN
1792         WITH CONGENERES[1] DO GETSETPOSNS(POSNS1, ROWNUM, BOND);
1793         LISTPOSNS(PRNTPOSNS^.COMBMEMS, POSNS1, [], []);
1794         IF REVERSIBLE THEN
1795             LISTPOSNS(PRNTPOSNS^.COMBMEMS, [], POSNS1, [])
1796         END
1797     END;
1798     FILIAL      : PROGERROR(8); {substituent node with filial congener}

```

```

1799          PARENTAL :
1800          CASE CONGENERS[2].RELATIONSHIP OF
1801          FRATERNAL : BEGIN
1802                  WITH CONGENERS[2] DO GETSETPOSNS(POSNS2, ROWNUM, BOND);
1803                  LISTPOSNS(PRNTPOSNS^.COMBMEMS, [0], POSNS2, []);
1804                  IF REVERSIBLE THEN
1805                      LISTPOSNS(PRNTPOSNS^.COMBMEMS, POSNS2, [0], [])
1806                  END;
1807          FILIAL    : PROGERROR(9); {substituent node with filial congner }
1808          PARENTAL : BEGIN
1809                  NEW(PRNTPOSNS^.COMBMEMS);
1810                  ECTRSIZE := ECTRSIZE + 6;
1811                  WITH PRNTPOSNS^.COMBMEMS^ DO
1812                      BEGIN
1813                          FIRST := 0;
1814                          SECOND := 0;
1815                          NEXT := NIL
1816                      END;
1817          END
1818          END
1819          END
1820
1821      END;
1822      {.....}
1823
1824
1825
1826      PROCEDURE DECLAREMULT(MULTIP    : MULTIPLIER;
1827                          MULTSUBST : SUBSTITUENT);
1828
1829      { Adds an entry to MDECLARATIONTABLE for MULTIP.
1830        Called by Body of PROCESSCT}
1831
1832      VAR PMPTR : PMDECLIST;
1833
1834      BEGIN
1835      NEW(PMPTR);
1836      WITH PMPTR^ DO
1837          BEGIN
1838              SUBSTDECN := RDECLARATIONTABLE[MULTSUBST];

```

APPENDIX 3:

GENSAL INTERPRETER

```

1839         NEXT := MDECLARATIONTABLE[MULTIP]
1840     END;
1841     MDECLARATIONTABLE[MULTIP] := PMPTR;
1842     DECLMULT := DECLMULT + [MULTIP]
1843     END;
1844
1845
1846
1847     BEGIN (* Body of Procedure PROCESSCT *)
1848     REPEAT
1849         IF INTERACTIVE THEN READFELDMANN;
1850         APICCOUNT := 0;
1851         REJECTED := FALSE;
1852         FOR ROWNO := 1 TO MAXCT DO IF ROWNO > NUMOFNODES
1853             THEN CT[ROWNO] := NIL
1854             ELSE CASE NODENATURE(ROWNO) OF
1855                 APICLABEL      : BEGIN
1856                     CT[ROWNO] := NIL;
1857                     IF FELDCT[ROWNO].AR[2] <> 0 THEN REJECT(58, ROWNO);
1858                     IF FELDCT[ROWNO].MULT <> 0 THEN REJECT(57, ROWNO);
1859                     IF APICCOUNT = 2
1860                         THEN REJECT(59, 0)
1861                         ELSE APICCOUNT := APICCOUNT + 1
1862                 END;
1863                 VARPOSNLABEL,
1864                 EXPHYDROGEN  : BEGIN
1865                     CT[ROWNO] := NIL;
1866                     IF FELDCT[ROWNO].AR[2] <> 0 THEN REJECT(58, ROWNO);
1867                     IF FELDCT[ROWNO].MULT <> 0 THEN REJECT(57, ROWNO)
1868                 END;
1869                 ATOMIC        : BEGIN
1870                     IF FELDCT[ROWNO].MULT <> 0 THEN REJECT(57, ROWNO);
1871                     NEW(CT[ROWNO], TRUE);
1872                     ECTRSIZE := ECTRSIZE + 30;
1873                     WITH CT[ROWNO]^ DO
1874                         BEGIN
1875                             ATOMICROW := TRUE;
1876                             FOR M := 1 TO 2 DO ATOM[M] := FELDCT[ROWNO].CHEM[M];
1877                             CHARGE := FELDCT[ROWNO].CHGE;
1878                             READCONGENERS(CONGENERS, HYDROGENS, ATOMICROW, ROWNO);

```

```

1879         IF HYDROGENS=0 THEN HNUMBER(ROWNO)
1880         END
1881         END;
1882         SUBSTNODE : BEGIN
1883             NEW(CT[ROWNO], FALSE);
1884             ECTRSIZE := ECTRSIZE + 33;
1885             WITH CT[ROWNO]^ DO
1886                 BEGIN
1887                     ATOMICROW := FALSE;
1888                     NAME := SUBSTNAME(FELDCT[ROWNO].CHEM);
1889                     CHARGE := FELDCT[ROWNO].CHGE;
1890                     VALUES := NIL;
1891                     READCONGENERS(CONGENERS, HYDROGENS, ATOMICROW, ROWNO);
1892                     IF NAME IN DECLSUBS
1893                         THEN WITH RDECLARATIONTABLE[NAME]^ .CONN BONDS DO
1894                             IF (CONNECTIONS = NUMOFCONNS(CONGENERS, HYDROGENS, ROWNO))
1895                                 OR (CONNECTIONS = NOTSET)
1896                                     THEN {accords with previous declaration}
1897                                     ELSE REJECT(53, ROWNO)
1898                             ELSE CHECKEARLIERDEFN(NAME, NUMOFCONNS(CONGENERS, HYDROGENS, ROWNO), ROWNO)
1899                         END
1900                     END
1901             END;
1902         IF REJECTED
1903             THEN IF INTERACTIVE
1904                 THEN BEGIN
1905                     FOR ROWNO := 1 TO NUMOFNODES DO IF CT[ROWNO] <> NIL THEN
1906                         IF CT[ROWNO]^ .ATOMICROW
1907                             THEN BEGIN
1908                                 DISPOSE(CT[ROWNO], TRUE);
1909                                 ECTRSIZE := ECTRSIZE - 30
1910                             END
1911                         ELSE BEGIN
1912                             DISPOSE(CT[ROWNO], FALSE);
1913                             ECTRSIZE := ECTRSIZE - 33
1914                         END;
1915                     FELDMODE := OLDDIAGRAM;
1916                     FELDMN(FELDMODE, FELDFIL);
1917                     IF FELDMODE = OLDDIAGRAM THEN FAILURE(41, 0, ' ')
1918                 END

```

APPENDIX 3:

GENSAL INTERPRETER

```

1919         ELSE FAILURE(41, 0, '  ')
1920 UNTIL NOT REJECTED;
1921
1922 FOR ROWNO := 1 TO NUMOFNODES DO IF CT[ROWNO] <> NIL THEN WITH CT[ROWNO]^ DO
1923   IF NOT ATOMICROW THEN
1924     BEGIN
1925       GETPOSNS(CONGENERS, CONNBONDS, PRNTPOSNS);
1926       DECLARESUBST(CT[ROWNO]^ .NAME,
1927                   PSADDRESS,
1928                   NIL,
1929                   CONNBONDS,
1930                   PRNTPOSNS);
1931       IF FELDCT[ROWNO].MULT <> 0 THEN DECLAREMULT(FELDCT[ROWNO].MULT, CT[ROWNO]^ .NAME);
1932       IF INDEPENDENT(ROWNO)
1933         THEN WITH CT[ROWNO]^ DO ENTERCOMBIN(NAME, VALUES)
1934         ELSE BEGIN
1935           ENTERCOMBIN(CT[ROWNO]^ .NAME, PSADDRESS^.CHILDGATE);
1936           DISPOSE(CT[ROWNO], FALSE);
1937           ECTRSIZE := ECTRSIZE - 33
1938         END
1939     END
1940 END;
1941
1942 -----}
1943
1944
1945
1946
1947 {-----}
1948 PROCEDURE READSD(VAR PSADDRESS : PTRPSTYPE;
1949                 INTERACTIVE   : BOOLEAN);
1950
1951 { Sets up a SPECIFIC partial structure in PSADDRESS, uses SPLITLINE and
1952   DIVIDELINE to handle Gensal lines containing tokens after the SD, and calls
1953   PROCESSCT to reformat the connection table. If INTERACTIVE is TRUE then
1954   FELDMN and READFELDMANN are used to produce FELDCT and FELDBD. Otherwise
1955   they are derived by DECODECT.
1956   Called by ALTNVLIST\ELEMENT
1957   Body of INTERPRET}
1958

```

```

1959 VAR OLDLINE      : LINELIST;    { GENSAL source line from which READSD was called }
1960     LINECONTINUED : BOOLEAN;    { Indicates more GENSAL on line }
1961
1962
1963
1964 FUNCTION SPLITLINE : BOOLEAN;
1965
1966     { TRUE if there is any non-space character beyond the current position (N) in
1967     CURRENTLINE, which is space-filled from the current position, the original
1968     version being saved in OLDLINE }
1969
1970     VAR M : 0..MAXLENGTH;
1971
1972     BEGIN
1973     OLDLINE := CURRENTLINE^;
1974     SPLITLINE := FALSE;
1975     FOR M := N TO MAXLENGTH DO
1976     IF CURRENTLINE^.LINE[M] <> ' ' THEN
1977     BEGIN
1978     SPLITLINE := TRUE;
1979     CURRENTLINE^.LINE[M] := ' '
1980     END
1981     END;
1982
1983
1984
1985 PROCEDURE DIVIDELINE(VAR CURRENTLINE : PLINELIST);
1986
1987     { Places the second half of OLDLINE.LINE in a new location in the
1988     linked list of lines }
1989
1990     VAR M : 1..MAXLENGTH;
1991
1992     BEGIN
1993     FOR M := 1 TO (N-1) DO OLDLINE.LINE[M] := ' ';
1994     OLDLINE.LAST := CURRENTLINE;
1995     OLDLINE.NEXT := CURRENTLINE^.NEXT;
1996     NEW(CURRENTLINE^.NEXT);
1997     CURRENTLINE^.NEXT^ := OLDLINE;
1998     CURRENTLINE := CURRENTLINE^.NEXT

```

```

1999      END;
2000
2001
2002
2003      BEGIN { body of procedure READSD }
2004      NEW(PSADDRESS, SPECIFIC);
2005      ECTRSIZE := ECTRSIZE + 70;
2006      WITH PSADDRESS^ DO
2007          BEGIN
2008              PSVARIETY := SPECIFIC;
2009              VISITED := FALSE;
2010              CHILDGATE := NIL;
2011              PARENTGATE := NIL
2012          END;
2013      LINECONTINUED := SPLITLINE;
2014      IF INTERACTIVE
2015          THEN BEGIN
2016              FELDMODE := NEWDIAGRAM;
2017              WRITELN;
2018              WRITELN('FELDMANN graphics system for structure diagram input and display:');
2019              FELDMN(FELDMODE,FELDFIL)
2020          END
2021          ELSE BEGIN
2022              CURRENTLINE := CURRENTLINE^.NEXT;
2023              DECODECT(CURRENTLINE,TRUE);
2024              IF LINECONTINUED THEN DIVIDELINE(CURRENTLINE)
2025                  ELSE N := MAXLENGTH
2026          END;
2027      PROCESSCT(PSADDRESS^.CT, INTERACTIVE, PSADDRESS);
2028      IF INTERACTIVE THEN
2029          BEGIN
2030              NEW(CURRENTLINE^.NEXT);
2031              CURRENTLINE^.NEXT^.LAST := CURRENTLINE;
2032              CURRENTLINE := CURRENTLINE^.NEXT;
2033              CURRENTLINE^.NEXT := NIL;
2034              ENCODECT(CURRENTLINE);
2035              IF LINECONTINUED THEN DIVIDELINE(CURRENTLINE)
2036                  ELSE N := MAXLENGTH
2037          END
2038      END;

```

```

2039          { of procedure READSD
2040 -----}
2041
2042
2043
2044 FUNCTION CHECKDELIM (VALIDELIMS : DELIMSET) : DELIMTYPE;
2045
2046 BEGIN
2047     IF TOKEN.NATURE = DELIMITER
2048     THEN IF TOKEN.DELIMVAL IN VALIDELIMS
2049           THEN CHECKDELIM := TOKEN.DELIMVAL
2050           ELSE CHECKDELIM := INVALIDTOKEN
2051     ELSE CHECKDELIM :=INVALIDTOKEN
2052 END;
2053
2054
2055
2056 {-----}
2057
2058           P R O C E D U R E   I N T E G E R R A N G E
2059 -----}
2060
2061 PROCEDURE INTEGERRANGE (VAR RANGEVALUES : INTRECORD;
2062                        LIMITRANGE      : INTRECORD;
2063                        ERRORCODE       : INTEGER);
2064
2065 { Carries out syntactic and semantic checking on integer ranges. LIMITRANGE is
2066   the range of values that all values in RANGEVALUES must fall, and is used for
2067   the semantic checking (functions INCREASING, WITHINLIMITS and ALLWITHINLIMITS).
2068   ERRORCODE is the relevant error code for passing to procedure ERROR.
2069   Called by GROUPRANGE
2070   Selector}
2071
2072
2073 VAR PTR : PDOUBLIST;
2074
2075
2076 FUNCTION WITHINLIMITS(TESTVALUE : INTEGER) : BOOLEAN;
2077
2078 { Returns TRUE is TESTVALUE is in the range covered by LIMITRANGE

```

```

2079         Called by ALLWITHINLIMITS
2080             RANGEFRAGMENT}
2081
2082     VAR PTR : PDOUBLIST;
2083
2084     BEGIN
2085     PTR := LIMITRANGE.SUBRANGES;
2086     WITHINLIMITS := FALSE;
2087     IF (TESTVALUE >= LIMITRANGE.TOPRANGE) AND (LIMITRANGE.TOPRANGE >= 0)
2088     THEN WITHINLIMITS := TRUE
2089     ELSE WHILE PTR <> NIL DO
2090         IF (TESTVALUE > PTR^.SECOND)
2091         THEN PTR := NIL
2092         ELSE IF TESTVALUE < PTR^.FIRST
2093         THEN PTR := PTR^.NEXT
2094         ELSE BEGIN
2095             WITHINLIMITS := TRUE;
2096             PTR := NIL
2097         END
2098     END;
2099
2100
2101
2102     FUNCTION INCREASING (TESTVALUE : INTEGER ) : BOOLEAN;
2103
2104     { Returns TRUE is TESTVALUE is larger than than the last integer in the range
2105     Called by RANGEFRAGMENT}
2106
2107     BEGIN
2108     IF RANGEVALUES.SUBRANGES = NIL
2109     THEN INCREASING := TRUE { This is the first integer in the range }
2110     ELSE INCREASING := TESTVALUE > RANGEVALUES.SUBRANGES^.SECOND
2111     END;
2112
2113
2114
2115     FUNCTION ALLWITHINLIMITS(LOWERBOUND,
2116                             UPPERBOUND : INTEGER): BOOLEAN;
2117
2118     { Returns TRUE is all the values between LOWERBOUND and UPPERBOUND inclusive

```

```

2119         are covered by LIMITRANGE.
2120         Called by RANGEFRAGMENT}
2121
2122     VAR VALID : BOOLEAN;
2123
2124     BEGIN
2125     VALID := TRUE;
2126     WHILE (LOWERBOUND <= UPPERBOUND) AND VALID DO
2127         IF WITHINLIMITS(LOWERBOUND)
2128             THEN LOWERBOUND := LOWERBOUND + 1
2129             ELSE BEGIN
2130                 ERROR(ERRORCODE, LOWERBOUND);
2131                 VALID := FALSE
2132             END;
2133     ALLWITHINLIMITS := VALID
2134     END;
2135
2136
2137
2138     PROCEDURE RANGEFRAGMENT;
2139
2140     { Carries out syntactic/semantic checking on a single range fragment. On
2141     entry to the procedure TOKEN is the token immediately before the first
2142     integer of the fragment. On exit, TOKEN is a comma or integer range
2143     terminating token.
2144     Called by Body of INTEGERRANGE}
2145
2146     VAR TERMINATORS : DELIMSET; { Tokens that terminate an integer range }
2147         VALID : BOOLEAN;
2148         FIRSTINTEGER : INTEGER; { The first integer in N1-N2 type ranges }
2149         DELIMCHECK : DELIMTYPE;
2150
2151     BEGIN
2152     TERMINATORS := [GCLOSANG, GRSQUARE, GEQUALS, GSEQ, GDEQ, GDOLEQ, GHASHEQ];
2153     NEXTTOKEN;
2154     REPEAT
2155         VALID := FALSE;
2156         WHILE TOKEN.NATURE <> INTEGRAL DO ERROR(23,0);
2157         IF NOT INCREASING(TOKEN.INTEGVAL)
2158             THEN ERROR(27,0)

```

```

2159         ELSE IF WITHINLIMITS(TOKEN.INTEGVAL)
2160             THEN VALID := TRUE
2161             ELSE ERROR(ERRORCODE,TOKEN.INTEGVAL);
2162 UNTIL VALID;
2163 FIRSTINTEGER := TOKEN.INTEGVAL;
2164 NEXTTOKEN;
2165 REPEAT
2166     VALID := FALSE;
2167     DELIMCHECK := CHECKDELIM([GCOMMA,GHYPHEN]+TERMINATORS);
2168     IF DELIMCHECK=INVALIDTOKEN THEN ERROR(24,0)
2169         ELSE VALID := TRUE
2170 UNTIL VALID;
2171 IF DELIMCHECK <> GHYPHEN
2172     THEN ADDINTS(RANGEVALUES.SUBRANGES, FIRSTINTEGER, FIRSTINTEGER)
2173     ELSE BEGIN
2174         NEXTTOKEN;
2175         REPEAT
2176             VALID := FALSE;
2177             WHILE (TOKEN.NATURE <> INTEGRAL) AND (CHECKDELIM(TERMINATORS) = INVALIDTOKEN)
2178                 DO ERROR(24,0);
2179             IF TOKEN.NATURE = INTEGRAL
2180                 THEN IF TOKEN.INTEGVAL < FIRSTINTEGER
2181                     THEN ERROR(27,0)
2182                     ELSE IF ALLWITHINLIMITS(FIRSTINTEGER, TOKEN.INTEGVAL)
2183                         THEN BEGIN
2184                             VALID := TRUE;
2185                             ADDINTS(RANGEVALUES.SUBRANGES, FIRSTINTEGER, TOKEN.INTEGVAL);
2186                             NEXTTOKEN;
2187                             WHILE CHECKDELIM([GCOMMA] + TERMINATORS)= INVALIDTOKEN DO ERROR(24, 0)
2188                             END
2189                         ELSE
2190                             ELSE BEGIN
2191                                 IF LIMITRANGE.TOPRANGE = NOTSET
2192                                     THEN IF LIMITRANGE.SUBRANGES=NIL
2193                                         THEN ERROR(ERRORCODE,0)
2194                                         ELSE ERROR(ERRORCODE, LIMITRANGE.SUBRANGES^.SECOND + 1)
2195                                     ELSE VALID := ALLWITHINLIMITS(FIRSTINTEGER, LIMITRANGE.TOPRANGE);
2196                                 IF VALID THEN RANGEVALUES.TOPRANGE := FIRSTINTEGER
2197                             END
2198 UNTIL VALID

```

```

2199         END
2200     END;
2201
2202
2203     BEGIN { Body of INTEGERRANGE }
2204     RANGEVALUES.SUBRANGES := NIL;
2205     RANGEVALUES.TOPRANGE := NOTSET;
2206     REPEAT RANGEFRAGMENT
2207     UNTIL TOKEN.DELIMVAL <> GCOMMA;
2208     PTR := RANGEVALUES.SUBRANGES;
2209     WHILE PTR <> NIL DO
2210         BEGIN
2211             ECTRSIZE := ECTRSIZE + 6;
2212             PTR := PTR^.NEXT
2213         END
2214     END;
2215         { of Procedure INTEGERRANGE
2216
2217     -----}
2218
2219
2220
2221
2222     PROCEDURE SETINTS (VAR RANGE : INTRECORD;
2223                       ONESET : INTEGSET);
2224
2225     { Takes a set of integers, and converts them to integer range format. If
2226     MAXVARS is a member of the set, then TOPRANGE is set to the member
2227     of the set above the highest absent member.
2228     Called from GROUPRANGE
2229     SELECTOR}
2230
2231     VAR N : NOTSET..MAXVARS;
2232
2233     BEGIN
2234     WITH RANGE DO
2235         BEGIN
2236             IF MAXVARS IN ONESET
2237             THEN BEGIN
2238                 N := MAXVARS;

```

```

2239         WHILE N IN ONESET DO N := N-1;
2240         TOPRANGE := N + 1;
2241         ONESET := ONESET - [TOPRANGE..MAXVARS]
2242     END
2243     ELSE TOPRANGE := NOTSET;
2244     SUBRANGES := NIL;
2245     FOR N := 0 TO MAXVARS DO
2246         IF N IN ONESET THEN
2247             ADDINTS(RANGE.SUBRANGES, N,N);
2248     END
2249 END;
2250
2251
2252
2253 PROCEDURE INTSET(VAR ONESET : INTEGSET;
2254                 RANGE      : INTRECORD);
2255
2256 { Converts an integer range into a set of integers, and DESTROYs the SUBRANGES
2257   of the integer range.
2258   Called from ALTNVLIST\ELEMENT\PARAMETERLIST}
2259
2260 VAR PTR : PDOUBLIST;
2261     M   : INTEGER;
2262
2263 BEGIN
2264 WITH RANGE DO
2265     BEGIN
2266     IF TOPRANGE=NOTSET
2267     THEN ONESET := []
2268     ELSE ONESET := [TOPRANGE..MAXVARS];
2269     PTR := SUBRANGES;
2270     WHILE PTR <> NIL DO WITH PTR^ DO
2271     BEGIN
2272     FOR M := FIRST TO SECOND DO ONESET := ONESET + [M];
2273     PTR := NEXT
2274     END;
2275     REDUCEECTR(SUBRANGES);
2276     DESTROY(SUBRANGES)
2277     END
2278 END;

```

```

2279
2280
2281
2282 PROCEDURE GROUPRANGE (VAR MEMBERS : INTEGSET;
2283                       LIMITSET   : INTEGSET;
2284                       ERRORCODE  : INTEGER);
2285
2286 { Converts LIMITSET into a INTRECORD format, and uses this as the Limitrange
2287   for a call to INTEGERRANGE. The RANGE that this returns is converted back
2288   to a set (MEMBERS).
2289   Called by ASIGNMENTSTMNT\SUBSTGROUP
2290           ASSIGNMENTSTMNT\MULTASSIGNMENT
2291           ALTNVLIST\POSITIONSET}
2292
2293 VAR RANGE,
2294     LIMITRANGE : INTRECORD;
2295     PTR        : PDOUBLIST;
2296     VAL        : 0..MAXVARS;
2297
2298 BEGIN
2299 MEMBERS := [];
2300 SETINTS(LIMITRANGE,LIMITSET);
2301 INTEGERRANGE(RANGE,LIMITRANGE,ERRORCODE);
2302 PTR := RANGE.SUBRANGES;
2303 WHILE PTR <> NIL DO WITH PTR^ DO
2304   BEGIN
2305     FOR VAL := FIRST TO SECOND DO
2306       MEMBERS := MEMBERS + [VAL];
2307     PTR := NEXT
2308   END;
2309 REDUCEECTR(LIMITRANGE.SUBRANGES);
2310 DESTROY(LIMITRANGE.SUBRANGES);
2311 DESTROY(RANGE.SUBRANGES)
2312 END;
2313
2314
2315
2316 PROCEDURE CHECKVALIDINT (LIMITSET : INTEGSET;
2317                          ERRORCODE : INTEGER );
2318

```

```

2319 { Checks that the current TOKEN is an integer within LIMITSET, and obtains
2320 further tokens from the input stream if it is not.
2321 Called by ASSIGNMENTSTMNT\SUBSTGROUP\SUBSTCOMBINATION
2322 ALTNVLIST\POSITIONSET\POSNCOMBINATION
2323 ALTNVLIST\POSITIONSET}
2324
2325 VAR VALID : BOOLEAN;
2326
2327 BEGIN
2328 VALID := FALSE;
2329 REPEAT
2330 WHILE TOKEN.NATURE <> INTEGRAL DO ERROR(23,0);
2331 IF TOKEN.INTEGVAL > MAXVARS
2332 THEN ERROR(ERRORCODE, TOKEN.INTEGVAL)
2333 ELSE IF TOKEN.INTEGVAL IN LIMITSET
2334 THEN VALID := TRUE
2335 ELSE ERROR(ERRORCODE, TOKEN.INTEGVAL)
2336 UNTIL VALID
2337 END;
2338
2339
2340
2341 PROCEDURE SELECTOR(VAR VALUERANGE : INTRECORD;
2342 LIMITSET : INTEGSET;
2343 ERRORCODE : INTEGER);
2344
2345 { Analyses a Gensal selector, returning the values in VALUERANGE. Limited by
2346 LIMITSET. ERRORCODE is passed to INTEGERRANGE.
2347 Called from ALTNVLIST\ELEMENT\PARAMETERLIST\USERPARAMETER
2348 ALTNVLIST\ELEMENT\PARAMETERLIST
2349 ALTNVLIST\ELEMENT
2350 ASSIGNMENTSTMNT
2351 ASSIGNMENTSTMNT\MULTASSIGNMENT}
2352
2353 VAR LIMITRANGE : INTRECORD;
2354
2355 BEGIN
2356 WHILE CHECKDELIM([GOPENANG])=INVALIDTOKEN DO ERROR(21,0);
2357 SETINTS(LIMITRANGE, LIMITSET);
2358 INTEGERRANGE(VALUERANGE, LIMITRANGE, ERRORCODE);

```

```

2359 DESTROY (LIMITRANGE.SUBRANGES);
2360 WHILE CHECKDELIM([GCLOSANG])=INVALIDTOKEN DO ERROR(22,0)
2361 END;
2362
2363
2364
2365 {*****}
2366
2367           P R O C E D U R E   A L T N V L I S T
2368
2369 *****}
2370
2371 PROCEDURE ALTNVLIST(PARENTPSLIST      : PPSLIST;
2372                   OPTIONALSUB        : BOOLEAN);
2373
2374 { Processes alternatives separated by / delimiters.
2375   Called by ASSIGNMENTSTMNT\SUBSTASSIGNMENT
2376     ALTNVLIST\ELEMENT (recursively)
2377     ALTNVLIST\ELEMENT\TRANSLATENOMEN (recursively) }
2378
2379 TYPE PPALTBARS = ^PALTBAR;
2380     PALTBAR   = RECORD
2381         PARSTRUCT   : PTRPSTYPE;
2382         ALTBAR      : PALTERNLIST;
2383         CONNBONDS   : TCONNBONDS;
2384         PRNTPOSNS   : PTGROUPMEMS;
2385         COPYCHILDPS : BOOLEAN;
2386         NEXT        : PPALTBARS
2387     END;
2388
2389 VAR PARALTLIST,
2390     WRITEPTR      : PPALTBARS;
2391     NEWALTERNATIVE : PALTERNLIST;
2392     READPTR       : PPSLIST;
2393
2394
2395
2396 PROCEDURE UPDATEPARALTCONNS(PARALTLIST : PPALTBARS);
2397
2398 { Copies the CONNBONDS field of the first item in PARALTLIST into all the

```

```

2399     other items in the list.
2400     Called by ELEMENT\VALIDSUBST
2401     ELEMENT}
2402
2403     VAR NEWCONNBONDS : TCONNBONDS;
2404
2405     BEGIN
2406     NEWCONNBONDS := PARALTLIST^.CONNBONDS;
2407     REPEAT
2408     PARALTLIST^.CONNBONDS := NEWCONNBONDS;
2409     PARALTLIST := PARALTLIST^.NEXT
2410     UNTIL PARALTLIST = NIL
2411     END;
2412
2413
2414
2415     {-----}
2416     PROCEDURE POSITIONSET(VAR SETMEMS      : TGROUPMEMS;
2417     AVAILABLEPOSITIONS : TGROUPMEMS;
2418     CONNECTIVITY      : TCONNS;
2419     ERRORCODE         : INTEGER);
2420
2421     { Analyses a position set.
2422     Called from MODIFYCHILDPOSITIONS
2423     ELEMENT}
2424
2425     VAR AVAILFIRST : INTEGSET;
2426
2427
2428
2429     PROCEDURE FINDFIRST(VAR POSNSET : INTEGSET;
2430     POSNLIST      : PDOUBLIST);
2431
2432     { Returns a set consisting of the FIRST fields of all the items in POSNLIST. }
2433
2434     BEGIN
2435     POSNSET := [];
2436     WHILE POSNLIST <> NIL DO WITH POSNLIST^ DO
2437     BEGIN
2438     POSNSET := POSNSET + [FIRST];

```

```

2439         POSNLIST := NEXT
2440     END
2441 END;
2442
2443
2444
2445     PROCEDURE FINDSECOND(VAR POSNSET : INTEGSET;
2446                          POSNLIST   : PDOUBLIST;
2447                          FIRSTPOSN  : ATOMNUMBER);
2448
2449     { Returns a set consisting of the SECOND fields of the items in POSNLIST
2450       that have FIRSTPOSN as FIRST field. }
2451
2452     BEGIN
2453     POSNSET := [];
2454     WHILE POSNLIST <> NIL DO WITH POSNLIST^ DO
2455     BEGIN
2456         IF FIRST = FIRSTPOSN THEN POSNSET := POSNSET + [SECOND];
2457         POSNLIST := NEXT
2458     END
2459     END;
2460
2461
2462
2463     PROCEDURE POSNCOMBINATION(AVAILPOSNS : INTEGSET;
2464                              VAR COMBMEMS : PDOUBLIST);
2465
2466     { Analyses a position combination, checking the validity of each position,
2467       and inserting it into the front of the list headed by COMBMEMS. }
2468
2469     VAR POSNPAIR : PDOUBLIST;
2470
2471     BEGIN
2472     NEW(POSNPAIR);
2473     ECTRSIZE := ECTRSIZE + 6;
2474     NEXTTOKEN;
2475     CHECKVALIDINT(AVAILPOSNS, ERRORCODE);
2476     POSNPAIR^.FIRST := TOKEN.INTEGVAL;
2477     NEXTTOKEN;
2478     WHILE CHECKDELIM([GSLASH]) <> GSLASH DO ERROR(33,0);

```

```

2479 IF AVAILABLEPOSITIONS.COMBINED
2480 THEN FINDSECOND(AVAILPOSNS, AVAILABLEPOSITIONS.COMBMEMS, POSNPAIR^.FIRST)
2481 ELSE {leave AVAILPOSNS the same};
2482 NEXTTOKEN;
2483 CHECKVALIDINT(AVAILPOSNS, ERRORCODE);
2484 POSNPAIR^.SECOND := TOKEN.INTEGVAL;
2485 POSNPAIR^.NEXT := COMBMEMS;
2486 COMBMEMS := POSNPAIR;
2487 NEXTTOKEN;
2488 WHILE CHECKDELIM([GCOMMA, GRSQUARE]) = INVALIDTOKEN DO ERROR(24,0)
2489 END;
2490
2491
2492 BEGIN { Body of Procedure POSITIONSET }
2493 IF AVAILABLEPOSITIONS.COMBINED
2494 THEN FINDFIRST(AVAILFIRST, AVAILABLEPOSITIONS.COMBMEMS)
2495 ELSE AVAILFIRST := AVAILABLEPOSITIONS.MEMBERS;
2496 LOOKAHEAD;
2497 CHECKVALIDINT(AVAILFIRST, ERRORCODE);
2498 LOOKAHEAD;
2500 CASE CONNECTIVITY OF
2501 NOTSET : WHILE CHECKDELIM([GSLASH, GCOMMA, GHYPHEN, GRSQUARE]) = INVALIDTOKEN DO ERROR(24,0);
2502 0 : PROGERROR(10); {attempting to process position set for unconnected substituent}
2503 1 : WHILE CHECKDELIM([GCOMMA, GHYPHEN, GRSQUARE]) = INVALIDTOKEN DO
2504 IF CHECKDELIM([GSLASH])=GSLASH THEN ERROR(34,0)
2505 ELSE ERROR(24,0);
2506 2 : WHILE CHECKDELIM([GSLASH]) <> GSLASH DO ERROR(33,0)
2507 END;
2508 IF TOKEN.DELIMVAL = GSLASH
2509 THEN BEGIN
2510 SETMEMS.COMBINED := TRUE;
2511 SETMEMS.COMBMEMS := NIL;
2512 NEXTTOKEN;
2513 REPEAT POSNCOMBINATION(AVAILFIRST, SETMEMS.COMBMEMS)
2514 UNTIL CHECKDELIM([GRSQUARE]) = GRSQUARE
2515 END
2516 ELSE BEGIN
2517 NEXTTOKEN;
2518 SETMEMS.COMBINED := FALSE;

```

```

2519         GROUPRANGE (SETMEMS.MEMBERS, AVAILFIRST, ERRORCODE)
2520     END;
2521 IF AVAILABLEPOSITIONS.COMBINED
2522     THEN BEGIN
2523         REDUCEECTR (AVAILABLEPOSITIONS.COMBMEMS);
2524         DESTROY (AVAILABLEPOSITIONS.COMBMEMS)
2525     END
2526 END;
2527     { of Procedure POSITIONSET
-----}
2528
2529
2530
2531
2532 FUNCTION COPYLIST (COMBMEMS : PDOUBLIST) : PDOUBLIST;
2533
2534 { Makes a reversed copy of COMBMEMS
2535   Called by ALTNVLIST\MODIFYCHILDPOSITIONS\TRACEDOWNGATE
2536   ALTNVLIST\ELEMENT\SETCOMBARS\CHECKCOMBPOSNS
2537   ALTNVLIST\ELEMENT\GETLIMITPOSITIONS}
2538
2539 VAR NEWLIST,
2540     NEWITEM : PDOUBLIST;
2541
2542 BEGIN
2543     NEWLIST := NIL;
2544     WHILE COMBMEMS <> NIL DO WITH COMBMEMS^ DO
2545         BEGIN
2546             NEW (NEWITEM);
2547             ECTRSIZE := ECTRSIZE + 6;
2548             NEWITEM^.FIRST := FIRST;
2549             NEWITEM^.SECOND := SECOND;
2550             NEWITEM^.NEXT := NEWLIST;
2551             NEWLIST := NEWITEM;
2552             COMBMEMS := NEXT
2553         END;
2554     COPYLIST := NEWLIST
2555 END;
2556
2557
2558

```

```

2559 PROCEDURE REDUCE(VAR LIMITLIST : PDOUBLIST;
2560                   COMPSET      : TGROUPMEMS);
2561
2562 { Removes those items in LIMITLIST that do not appear in COMPSET
2563   Called by ALTNVLIST\MODIFYCHILDPOSITIONS\TRACEDOWNGATE
2564   ALTNVLIST\ELEMENT\TRANSLATENOMEN\MODIFYGATEPOSITIONS
2565   ALTNVLIST\ELEMENT\GETLIMITPOSITIONS}
2566
2567 VAR LISTPTR,
2568     LASTPTR,
2569     COMPPTR : PDOUBLIST;
2570     FOUND   : BOOLEAN;
2571
2572 BEGIN
2573 LISTPTR := LIMITLIST;
2574 LASTPTR := NIL;
2575 WHILE LISTPTR <> NIL DO
2576   BEGIN
2577     IF COMPSET.COMBINED
2578       THEN BEGIN
2579         FOUND := FALSE;
2580         COMPPTR := COMPSET.COMBMEMS;
2581         WHILE (COMPPTR <> NIL) AND NOT FOUND DO
2582           BEGIN
2583             FOUND := (COMPPTR^.FIRST=LISTPTR^.FIRST) AND (COMPPTR^.SECOND=LISTPTR^.SECOND);
2584             COMPPTR := COMPPTR^.NEXT
2585           END
2586         END
2587       ELSE FOUND := [LISTPTR^.FIRST, LISTPTR^.SECOND] <= COMPSET.MEMBERS;
2588
2589     IF FOUND THEN BEGIN
2590       LASTPTR := LISTPTR;
2591       LISTPTR := LASTPTR^.NEXT
2592     END
2593     ELSE IF LASTPTR = NIL
2594       THEN BEGIN
2595       LIMITLIST := LISTPTR^.NEXT;
2596       DISPOSE(LISTPTR);
2597       ECTRSIZE := ECTRSIZE - 6;
2598       LISTPTR := LIMITLIST

```

```
2599             END
2600         ELSE BEGIN
2601             LASTPTR^.NEXT := LISTPTR^.NEXT;
2602             DISPOSE(LISTPTR);
2603             ECTRSIZE := ECTRSIZE - 6;
2604             LISTPTR := LASTPTR^.NEXT
2605         END
2606     END
2607 END;
2608
2609
2610
2611 PROCEDURE CONCATENATETERMS(VAR GATEPS : PTRPSTYPE);
2612 { Sets up a PS of variety OTHER, and concatenates NOMENCLATURE tokens up to a
2613   maximum of TERMLENGTH chars into it.
2614   Called from ELEMENT}
2615
2616 VAR DELIMCHECK : DELIMTYPE;
2617     TERMEND     : BOOLEAN;
2618     M, M2       : 0..TERMLENGTH;
2619
2620
2621 BEGIN
2622 NEW(GATEPS, OTHER);
2623 ECTRSIZE := ECTRSIZE + 22;
2624 WITH GATEPS^ DO
2625     BEGIN
2626         PSVARIETY := OTHER;
2627         VISITED := FALSE;
2628         CHILDGATE := NIL;
2629         PARENTGATE := NIL;
2630         FOR M := 1 TO TERMLENGTH DO TERM[M] := ' '
2631     END;
2632 M := 0;
2633 NEXTTOKEN;
2634 REPEAT
2635     DELIMCHECK := CHECKDELIM([GPRIME]);
2636     IF DELIMCHECK=INVALIDTOKEN
2637     THEN IF TOKEN.NATURE <> NOMENCLATURE
2638     THEN ERROR(25,0)
```

```

2639         ELSE BEGIN
2640             TERMEND := M = TERMLENGTH;
2641             M2 := 0;
2642             WHILE NOT TERMEND DO
2643                 BEGIN
2644                     M := M + 1;
2645                     M2 := M2 + 1;
2646                     GATEPS^.TERM[M] := TOKEN.NOMENVAL[M2];
2647                     TERMEND := (M=TERMLENGTH) OR (TOKEN.NOMENVAL[M2]=' ')
2648                 END;
2649             NEXTTOKEN
2650         END
2651 UNTIL DELIMCHECK=GPRIME
2652 END;
2653
2654
2655
2656 FUNCTION RECORDHELD(TERM          : STRING32;
2657                     VAR ADDRESS : INTEGER) : BOOLEAN;
2658
2659 { Determines whether or not a record is held for TERM. Requests synonyms
2660   for the term (using TERMREAD) if initially unsuccessful.
2661   The search is abandoned if a record is found, or if TERMREAD returns
2662   FALSE.
2663   Called by ELEMENT\TRANSLATENOMEN}
2664
2665 VAR STILLLOOKING,
2666     SYNONYMREAD : BOOLEAN;
2667
2668 BEGIN
2669 SYNONYMREAD := FALSE;
2670 REPEAT
2671     STILLLOOKING := NORECORD(TERM, ADDRESS);
2672     RECORDHELD := NOT STILLLOOKING;
2673     IF STILLLOOKING
2674     THEN BEGIN
2675         WRITE('No record held for ');
2676         PRINTNOM(TERM);
2677         WRITELN(''.');
2678         STILLLOOKING := FALSE;

```

```

2679         CASE INPUTMODE OF
2680             TERMINAL      : BEGIN
2681                 WRITE('Enter synonym or <CR>: > ');
2682                 STILLLOOKING := TERMREAD(TERM);
2683                 SYNONYMREAD := STILLLOOKING
2684             END;
2685             STOREDGENSAL : ;
2686             INSERTTEXT  : WRITELN('(Term in inserted Gensal expression)')
2687         END
2688     END
2689 ELSE IF SYNONYMREAD
2690     THEN BEGIN
2691         WRITE('Record found for ');
2692         PRINTNOM(TERM);
2693         WRITELN('".');
2694     END
2695 UNTIL NOT STILLLOOKING
2696 END;

2699
2700 FUNCTION DEFNTABLEENTRY(PARSTRUCT : PTRPSTYPE) : BOOLEAN;
2701
2702 { Returns TRUE if PARSTRUCT is NIL or has no PARENTGATE (provided it is not
2703   INTERNALREP.CONSTANTPART). Since the parameter passed is the PARSTRUCT
2704   field of a PARALTLIST element, this indicates whether or not it is in the
2705   chain of structures pointed at by RDEFINITIONTABLE.
2706   Called by ELEMENT\SETCOMBARS
2707           ELEMENT\SUBSTASVALUE
2708           ADDFURTHERSUBTN}
2709
2710 BEGIN
2711 IF PARSTRUCT = NIL
2712 THEN DEFNTABLEENTRY := TRUE
2713 ELSE IF PARSTRUCT^.PARENTGATE = NIL
2714 THEN DEFNTABLEENTRY := PARSTRUCT <> INTERNALREP.CONSTANTPART
2715 ELSE DEFNTABLEENTRY := FALSE
2716 END;
2717
2718

```

```

2719
2720 PROCEDURE FINDPOSITIONS(PTRPS          : PTRPSTYPE;
2721                        VAR AVAILPOSNS  : INTEGSET;
2722                        BOND MAG        : TBOND MAG);
2723
2724 { Returns the positions in PTRPS^ which are substitutable by a bond of
2725  magnitude BOND MAG
2726  Called by GETCHILDPOSITIONS
2727          MODIFYCHILDPOSITIONS\GETCOMBPOSNS}
2728
2729 VAR ROWNO      : ATOMNUMBER;
2730
2731 BEGIN
2732 WITH PTRPS^ DO CASE PSVARIETY OF
2733   DUMMY,
2734   UNKNOWN,
2735   OTHER      : AVAILPOSNS := [1..MAXCT];
2736   GENERIC    : AVAILPOSNS := [1];
2737   SPECIFIC   : BEGIN
2738                 AVAILPOSNS := [];
2739                 FOR ROWNO := 1 TO MAXCT DO IF CT[ROWNO] <> NIL
2740                   THEN IF CT[ROWNO]^ .HYDROGENS >= BOND MAG
2741                     THEN AVAILPOSNS := AVAILPOSNS + [ROWNO]
2742   END
2743 END
2744 END;
2745
2746
2747
2748 {.....}
2749 PROCEDURE GETCHILDPOSITIONS(PTRPS          : PTRPSTYPE;
2750                            VAR CONNBONDS   : TCONNBONDS;
2751                            VAR CHILDPOSITIONS : TGROUPMEMS);
2752
2753 { Makes initial determination of CHILDPOSITIONS for PTRPS (which points to a
2754  Child PS), also modifying CONNBONDS as necessary. The CHILDPOSITIONS field
2755  may be further modified by a post substituent value position set.
2756  Called by ELEMENT\SETCOMBARS
2757          ELEMENT\SUBSTASVALUE}
2758

```

```
2759 VAR POSNA,
2760     POSNB      : ATOMNUMBER;
2761     POSNSETA,
2762     POSNSETB,
2763     POSNSETC   : INTEGSET;
2764     MAGSUM     : INTEGER;
2765     NUMMARKERS : TCONNS;
2766     BONDA,
2767     BONDB      : BONDORDER;
2768     FAILSTRING : STRING4;
2769
2770
2771
2772 FUNCTION BONDHECK(PARENTBOND,
2773                  CHILDBOND : BONDORDER) : BONDORDER;
2774
2775 { Checks compatability of the two bonds, ejecting user to the editor if the
2776   bonds are found to be incompatible. The global variable BONDMATCHARRAY is
2777   used to check the compatibility.
2778   Called from THISWAYROUND
2779   Body of GETCHILDPPOSITIONS}
2780
2781 VAR NEWBOND : BONDORDER;
2782     FAILDATA : STRING4;
2783
2784 BEGIN
2785 NEWBOND := BONDMATCHARRAY[PARENTBOND, CHILDBOND];
2786 IF NEWBOND = NOTSPECIFIED
2787 THEN BEGIN
2788     FAILDATA[1] := BONDSTRING[PARENTBOND,1];
2789     FAILDATA[2] := BONDSTRING[PARENTBOND,2];
2790     FAILDATA[3] := BONDSTRING[CHILDBOND, 1];
2791     FAILDATA[4] := BONDSTRING[CHILDBOND, 2];
2792     FAILURE(42, 0, FAILDATA)
2793     END
2794 ELSE BONDHECK := NEWBOND
2795 END;
```

```

2799 PROCEDURE GETMARKEDPOSNS(CT          : CTYPE;
2800                               VAR POSNA,
2801                               POSNB : ATOMNUMBER;
2802                               VAR BONDA,
2803                               BONDB : BONDORDER);
2804
2805 { Returns those positions in CT which have PARENTAL bonds, with their bond
2806   orders. On entry the parameters are NPTFIXED or NOTSPECIFIED. PROCESSCT
2807   will only have permitted a maximum of two marked positions.}
2808
2809 VAR ROWNO : ATOMNUMBER;
2810     CNGNR : 1..MAXCONGENERERS;
2811
2812 BEGIN
2813 FOR ROWNO := 1 TO MAXCT DO IF CT[ROWNO] <> NIL
2814   THEN FOR CNGNR := 1 TO MAXCONGENERERS DO
2815     WITH CT[ROWNO]^CONGENERERS[CNGNR] DO
2816       IF RELATIONSHIP = PARENTAL
2817         THEN IF POSNA = NOTFIXED
2818           THEN BEGIN
2819             POSNA := ROWNO;
2820             BONDA := BOND
2821           END
2822         ELSE BEGIN
2823           POSNB := ROWNO;
2824           BONDB := BOND
2825         END
2826     END;
2827
2828
2829 FUNCTION HYDROGENPS(PTRPS : PTRPSTYPE) : BOOLEAN;
2830
2831 { Returns TRUE if PTRPS represents hydrogen (i.e. has no non-hydrogen atoms).}
2832
2833 VAR ROWNO : ATOMNUMBER;
2834
2835 BEGIN
2836 IF PTRPS^.PSVARIETY = SPECIFIC
2837   THEN BEGIN

```

```

2839             HYDROGENPS := TRUE;
2840             FOR ROWNO := 1 TO MAXCT DO
2841                 IF PTRPS^.CT[ROWNO] <> NIL THEN HYDROGENPS := FALSE
2842             END
2843         ELSE HYDROGENPS := FALSE
2844     END;
2845
2846
2847
2848     FUNCTION THISWAYROUND(PARENTA,
2849                         CHILDA,
2850                         PARENTB,
2851                         CHILDB : BONDORDER) : BOOLEAN;
2852
2853     { Determines whether or not the bonds in a doubly-connected child need to be
2854       reversed for compatability with the parent. If either way will do, the way
2855       given is prefered unless the other way round matches identical (as opposed
2856       to merely compatible) bonds. }
2857
2858     VAR FITSTHISWAY,
2859         FITSOTHERWAY : BOOLEAN;
2860
2861     BEGIN
2862     FITSTHISWAY := (BONDMATCHARRAY[PARENTA,CHILDA] <> NOTSPECIFIED) AND (BONDMATCHARRAY[PARENTB,CHILDB] <> NOT
SPECIFIED);
2863     FITSOTHERWAY := (BONDMATCHARRAY[PARENTA,CHILDB] <> NOTSPECIFIED) AND (BONDMATCHARRAY[PARENTB,CHILDA] <> NO
SPECIFIED);
2864     IF FITSTHISWAY
2865     THEN IF FITSOTHERWAY
2866         THEN THISWAYROUND := NOT ((PARENTA = CHILDB) AND (PARENTB = CHILDA))
2867         ELSE THISWAYROUND := TRUE
2868     ELSE IF FITSOTHERWAY
2869         THEN THISWAYROUND := FALSE
2870         ELSE {Bond match failure - use BONDCHECK to give error message}
2871             IF BONDMATCHARRAY[PARENTA, CHILDA] = NOTSPECIFIED
2872             THEN PARENTA := BONDCHECK(PARENTA, CHILDA)
2873             ELSE PARENTB := BONDCHECK(PARENTB, CHILDB)
2874     END;
2875
2876

```

APPENDIX 3:

GENSAL INTERPRETER

```

2877
2878 PROCEDURE FINDNONAPICPOSNS(CPARAM      : INTRECORD;
2879                          VAR POSNSET : INTEGSET);
2880
2881 { Returns a position set containing the possible "right-hand end" terminal
2882   positions in a GENERIC PS, based on the possible values for the ATOMCOUNT
2883   parameter, passed as CPARAM. }
2884
2885 VAR PTR : PDOUBLIST;
2886
2887 BEGIN
2888   POSNSET := [];
2889   PTR := CPARAM.SUBRANGES;
2890   WHILE PTR <> NIL DO WITH PTR^ DO
2891     BEGIN
2892       IF SECOND <= MAXCT
2893         THEN POSNSET := POSNSET + [FIRST..SECOND]
2894         ELSE IF FIRST <= MAXCT
2895           THEN POSNSET := POSNSET + [FIRST..MAXCT];
2896       PTR := NEXT
2897     END;
2898   IF CPARAM.TOPRANGE <> NOTSET
2899     THEN POSNSET := POSNSET + [CPARAM.TOPRANGE..MAXCT]
2900   END;
2901
2902
2903
2904 BEGIN {Body of GETCHILDPOSITIONS}
2905   POSNA := NOTFIXED;
2906   POSNB := NOTFIXED;
2907   BONDA := NOTSPECIFIED;
2908   BONDB := NOTSPECIFIED;
2909   IF PTRPS^.PSVARIETY = SPECIFIC
2910     THEN BEGIN
2911       GETMARKEDPOSNS(PTRPS^.CT, POSNA, POSNB, BONDA, BONDB);
2912       NUMMARKERS := ORD(POSNA<>NOTFIXED) + ORD(POSNB<>NOTFIXED)
2913     END
2914   ELSE NUMMARKERS := 0;
2915   CASE CONNBONDS.CONNECTIONS OF
2916     NOTSET : CASE NUMMARKERS OF

```

```

2917           2 : BEGIN
2918             CONNBONDS.CONNECTIONS := 2;
2919             CONNBONDS.BONDA := BONDHECK(CONNBONDS.BONDA, NOTSPECIFIED);
2920             CONNBONDS.BONDB := BONDHECK(CONNBONDS.BONDB, NOTSPECIFIED);
2921             CHILDPOSITIONS.COMBINED := TRUE;
2922             CHILDPOSITIONS.COMBMEMS := NIL;
2923             LISTPOSNS(CHILDPOSITIONS.COMBMEMS, [POSNA], [POSNB], [POSNA, POSNB]);
2924             IF (BONDA=BONDB) OR (BONDMATCHARRAY[BONDA,BONDB] IN [ANY, CHAIN, RING])
2925               THEN LISTPOSNS(CHILDPOSITIONS.COMBMEMS, [POSNB], [POSNA], [])
2926             END;
2927           1 : BEGIN
2928             CONNBONDS.CONNECTIONS := 1;
2929             CONNBONDS.BOND := BONDA;
2930             CHILDPOSITIONS.COMBINED := FALSE;
2931             CHILDPOSITIONS.MEMBERS := [POSNA]
2932           END;
2933           0 : BEGIN
2934             CONNBONDS.CONNECTIONS := 1; {assumption}
2935             CONNBONDS.BOND := CHAISING; {assumption}
2936             CHILDPOSITIONS.COMBINED := FALSE;
2937             FINDPOSITIONS(PTRPS, CHILDPOSITIONS.MEMBERS, 1)
2938           END
2939         END;
2940       0 : BEGIN
2941         CHILDPOSITIONS.COMBINED := FALSE;
2942         CHILDPOSITIONS.MEMBERS := []
2943       END;
2944       1 : BEGIN
2945         CHILDPOSITIONS.COMBINED := FALSE;
2946         CASE NUMMARKERS OF
2947           2 : FAILURE(44, 0, ' ');
2948           1 : BEGIN
2949             CHILDPOSITIONS.MEMBERS := [POSNA];
2950             CONNBONDS.BOND := BONDHECK(CONNBONDS.BOND, BONDA);
2951           END;
2952         0 : BEGIN
2953           FINDPOSITIONS(PTRPS, CHILDPOSITIONS.MEMBERS, MAGNITUDE(CONNBONDS.BOND));
2954           CONNBONDS.BOND := BONDHECK(CONNBONDS.BOND, NOTSPECIFIED)
2955         END
2956       END

```

```

2957     END;
2958     2 : BEGIN
2959         CHILDPOSITIONS.COMBINED := TRUE;
2960         CHILDPOSITIONS.COMBMEMS := NIL;
2961         CASE NUMMARKERS OF
2962             0 : BEGIN
2963                 FINDPOSITIONS(PTRPS, POSNSETA, MAGNITUDE(CONNBONDS.BONDA));
2964                 IF PTRPS^.PSVARIETY = GENERIC
2965                     THEN FINDNONAPICPOSNS(PTRPS^.PARAMLIST[ATOMCOUNT], POSNSETB)
2966                     ELSE FINDPOSITIONS(PTRPS, POSNSETB, MAGNITUDE(CONNBONDS.BONDB));
2967                 WITH CONNBONDS DO MAGSUM := MAGNITUDE(BONDA)+MAGNITUDE(BONDB);
2968                 IF (MAGSUM <= 3) AND (POSNSETA * POSNSETB <> [])
2969                     THEN FINDPOSITIONS(PTRPS, POSNSETC, MAGSUM)
2970                     ELSE POSNSETC := [];
2971                 POSNSETC := POSNSETA * POSNSETB * POSNSETC;
2972                 LISTPOSNS(CHILDPOSITIONS.COMBMEMS, POSNSETA, POSNSETB, POSNSETC);
2973                 CONNBONDS.BONDA := BONDHECK(CONNBONDS.BONDA, BONDA);
2974                 CONNBONDS.BONDB := BONDHECK(CONNBONDS.BONDB, BONDB)
2975             END;
2976             1 : BEGIN
2977                 FINDPOSITIONS(PTRPS, POSNSETB, MAGNITUDE(CONNBONDS.BONDB));
2978                 POSNSETA := [POSNA];
2979                 WITH CONNBONDS DO MAGSUM := MAGNITUDE(BONDA)+MAGNITUDE(BONDB);
2980                 IF (MAGSUM <= 3) AND (POSNSETA * POSNSETB <> [])
2981                     THEN FINDPOSITIONS(PTRPS, POSNSETC, MAGSUM)
2982                     ELSE POSNSETC := [];
2983                 POSNSETC := POSNSETA * POSNSETB * POSNSETC;
2984                 IF THISWAYROUND(CONNBONDS.BONDA, BONDA, CONNBONDS.BONDB, BONDB)
2985                     THEN BEGIN
2986                         LISTPOSNS(CHILDPOSITIONS.COMBMEMS, POSNSETA, POSNSETB, POSNSETC);
2987                         IF BONDMATCHARRAY[CONNBONDS.BONDB, BONDA] <> NOTSPECIFIED
2988                             THEN LISTPOSNS(CHILDPOSITIONS.COMBMEMS, POSNSETB, POSNSETA, [])
2989                         CONNBONDS.BONDA := BONDHECK(CONNBONDS.BONDA, BONDA);
2990                         CONNBONDS.BONDB := BONDHECK(CONNBONDS.BONDB, BONDB)
2991                     END
2992                 ELSE BEGIN
2993                     LISTPOSNS(CHILDPOSITIONS.COMBMEMS, POSNSETB, POSNSETA, POSNSETC);
2994                     CONNBONDS.BONDA := BONDHECK(CONNBONDS.BONDA, BONDB);
2995
2996                     CONNBONDS.BONDB := BONDHECK(CONNBONDS.BONDB, BONDA)

```

```

2997                                     END
2998                                     END;
2999     2 : IF THISWAYROUND(CONNBONDS.BONDA, BONDA, CONNBONDS.BONDB, BONDB)
3000         THEN BEGIN
3001             LISTPOSNS(CHILDPOSITIONS.COMBMEMS, [POSNA], [POSNB], [POSNA, POSNB]);
3002             IF (BONDA=BONDB) OR (BONDMATCHARRAY[BONDA,BONDB] IN [ANY, CHAIN, RING])
3003                 THEN BEGIN
3004                     CONNBONDS.BONDA := BONDHECK(CONNBONDS.BONDA, BONDHECK(BONDA,BONDB));
3005                     CONNBONDS.BONDB := BONDHECK(CONNBONDS.BONDB, BONDHECK(BONDA,BONDB));
3006                     LISTPOSNS(CHILDPOSITIONS.COMBMEMS, [POSNB], [POSNA], [])
3007                 END
3008             ELSE BEGIN
3009                 CONNBONDS.BONDA := BONDHECK(CONNBONDS.BONDA, BONDA);
3010                 CONNBONDS.BONDB := BONDHECK(CONNBONDS.BONDB, BONDB)
3011             END
3012         END
3013     ELSE BEGIN
3014         LISTPOSNS(CHILDPOSITIONS.COMBMEMS, [POSNB], [POSNA], [POSNA,POSNB]);
3015         CONNBONDS.BONDA := BONDHECK(CONNBONDS.BONDA, BONDB);
3016         CONNBONDS.BONDB := BONDHECK(CONNBONDS.BONDB, BONDA)
3017     END
3018     END {CASE}
3019     END
3020     END; {CASE}
3021     IF CHILDPOSITIONS.COMBINED
3022     THEN IF CHILDPOSITIONS.COMBMEMS = NIL
3023         THEN BEGIN
3024             FAILSTRING[1] := BONDSTRING[CONNBONDS.BONDA, 1];
3025             FAILSTRING[2] := BONDSTRING[CONNBONDS.BONDA, 2];
3026             FAILSTRING[3] := BONDSTRING[CONNBONDS.BONDB, 1];
3027             FAILSTRING[4] := BONDSTRING[CONNBONDS.BONDB, 2];
3028             FAILURE(52, 0, FAILSTRING)
3029         END
3030     ELSE IF (CHILDPOSITIONS.MEMBERS = []) AND NOT HYDROGENPS(PTRPS)
3031         THEN BEGIN
3032             FAILSTRING[1] := BONDSTRING[CONNBONDS.BOND, 1];
3033             FAILSTRING[2] := BONDSTRING[CONNBONDS.BOND, 2];
3034             FAILSTRING[3] := ' ';
3035             FAILSTRING[4] := ' ';
3036             FAILURE(43, 0, FAILSTRING)

```

```

3037                                     END
3038     END;                               {of GETCHILDPOSITIONS
3039     .....}
3040
3041
3042
3043
3044     {-----}
3045     PROCEDURE MODIFYCHILDPOSITIONS(PARALTLIST : PPALTBARS);
3046
3047     { This procedure modifies the CHILDPOSITIONS fields of the bottom bars of the
3048       childgates, and also those of the parentgates, in accordance with the values
3049       given in the post substituent value position set.
3050       Called from ELEMENT}
3051
3052     TYPE PGBLIST = ^TBGLIST;
3053           TBGLIST = RECORD
3054               GBOTTOM : PCOMBINLIST;
3055               NEXT    : PGBLIST
3056           END;
3057
3058     VAR PTR,
3059         GATEBOTTOMS      : PGBLIST;
3060         PG                : PPARENTLIST;
3061         CHILDGATEPOSITIONS,
3062         LIMITPOSITIONS   : TGROUPMEMS;
3063         LIMITINITIALISED : BOOLEAN;
3064         CONNECTIVITY     : TCONNS;
3065
3066
3067     PROCEDURE GETCOMBPOSNS(CHILDPS          : PTRPSTYPE;
3068                           POSNSA          : INTEGSET;
3069                           VAR COMBAVAILPOSNS : TGROUPMEMS);
3070
3071     { Returns COMBAVAILPOSNS with a COMBINED position set of all possible position
3072       pairs in CHILDPS having members of POSNSA as their first member.}
3073
3074     VAR POSNSB,
3075         POSNSC : INTEGSET;
3076

```

```

3077 BEGIN
3078 FINDPOSITIONS(CHILDPS, POSNSB, 1);
3079 FINDPOSITIONS(CHILDPS, POSNSC, 2);
3080 COMBAVAILPOSNS.COMBINED := TRUE;
3081 COMBAVAILPOSNS.COMBMEMS := NIL;
3082 LISTPOSNS(COMBAVAILPOSNS.COMBMEMS, POSNSA, POSNSB, POSNSC)
3083 END;
3084
3085
3086
3087 PROCEDURE TRACEDOWNGATE(COMBINBAR : PCOMBINLIST;
3088                          CONNSFIXED : BOOLEAN);
3089
3090 { Traces down a child gate, adding the BOTTOMBARs to the GATEBOTTOMS list.
3091   LIMITPOSITIONS (in MODIFYCHILDPOSITIONS) is also initialised or updated
3092   appropriately. If CONNSFIXED is FALSE, then the connectivity of 1 recorded
3093   in COMBINBAR^.CONNBONDS is only an assumption, and could be modified by the
3094   position set about to be read. Therefore LIMITPOSITIONS must be COMBINED,
3095   GETCOMBPOSNS identifying all the possible second positions for the first
3096   positions identified by GETCHILDPOSITIONS. }
3097
3098 VAR NEWGB          : PGBLIST;
3099     ALTERNBAR      : PALTERNLIST;
3100     SUBCB          : PCOMBINLIST;
3101     COMBAVAILPOSNS : TGROUPMEMS;
3102
3103 BEGIN
3104 IF COMBINBAR^.BOTTOMBAR
3105 THEN BEGIN
3106     NEW(NEWGB);
3107     NEWGB^.NEXT := GATEBOTTOMS;
3108     NEWGB^.GBOTTOM := COMBINBAR;
3109     GATEBOTTOMS := NEWGB;
3110     WITH COMBINBAR^ DO IF LIMITINITIALISED
3111     THEN IF LIMITPOSITIONS.COMBINED
3112     THEN IF CHILDPOSITIONS.COMBINED
3113     THEN REDUCE(LIMITPOSITIONS.COMBMEMS, CHILDPOSITIONS)
3114     ELSE IF CONNSFIXED
3115     THEN PROGERROR(11) {mismatched combined fields}
3116     ELSE BEGIN

```

```

3117                                     GETCOMBPOSNS(CHILDPS, CHILDPOSITIONS.MEMBERS, COMBAVAILPOSNS);
3118                                     REDUCE(LIMITPOSITIONS.COMBMEMS, COMBAVAILPOSNS);
3119                                     REDUCEECTR(COMBAVAILPOSNS.COMBMEMS);
3120                                     DESTROY(COMBAVAILPOSNS.COMBMEMS)
3121                                     END
3122                                     ELSE IF CHILDPOSITIONS.COMBINED
3123                                         THEN PROGERROR(12) {mismatched combined fields}
3124                                         ELSE LIMITPOSITIONS.MEMBERS := LIMITPOSITIONS.MEMBERS * CHILDPOSITIONS.MEMBERS
3125     ELSE BEGIN
3126         IF CHILDPOSITIONS.COMBINED
3127             THEN BEGIN
3128                 LIMITPOSITIONS.COMBINED := TRUE;
3129                 LIMITPOSITIONS.COMBMEMS := COPYLIST(CHILDPOSITIONS.COMBMEMS)
3130             END
3131             ELSE IF CONNSFIXED
3132                 THEN LIMITPOSITIONS := CHILDPOSITIONS
3133                 ELSE GETCOMBPOSNS(CHILDPS, CHILDPOSITIONS.MEMBERS, LIMITPOSITIONS);
3134             LIMITINITIALISED := TRUE
3135         END
3136     END
3137 ELSE BEGIN
3138     ALTERNBAR := COMBINBAR^.ALTERNATIVES;
3139     WHILE ALTERNBAR <> NIL DO
3140     BEGIN
3141         SUBCB := ALTERNBAR^.COMBINATION;
3142         WHILE SUBCB <> NIL DO
3143         BEGIN
3144             TRACEDOWNGATE(ALTERNBAR^.COMBINATION, CONNSFIXED);
3145             SUBCB := SUBCB^.NEXT
3146         END;
3147         ALTERNBAR := ALTERNBAR^.NEXT
3148     END
3149 END
3150 END;
3151
3152
3153
3154 PROCEDURE ALTERCONNBONDS(VAR CONNBONDS : TCONNBONDS);
3155 { Adds a second bond to CONNBONDS, changing CONNECTIONS to 2 }
3156

```

APPENDIX 3:

GENSAL INTERPRETER

```

3157
3158     VAR NEWCONNBONDS : TCONNBONDS;
3159
3160     BEGIN
3161     WITH NEWCONNBONDS DO
3162         BEGIN
3163             CONNECTIONS := 2;
3164             BONDA := CONNBONDS.BOND;
3165             BONDB := CHAISING
3166         END;
3167     CONNBONDS := NEWCONNBONDS
3168     END;
3169
3170
3171
3172
3173     BEGIN {Body of MODIFYCHILDPOSITIONS}
3174     LIMITINITIALISED := FALSE;
3175     CONNECTIVITY := PARALTLIST^.CONNBONDS.CONNECTIONS;
3176     GATEBOTTOMS := NIL;
3177     WHILE PARALTLIST <> NIL DO WITH PARALTLIST^ DO
3178         BEGIN
3179             IF ALTBAR = NIL
3180                 THEN TRACEDOWNGATE(PARSTRUCT^.CHILDGATE, (CONNECTIVITY<>NOTSET))
3181                 ELSE TRACEDOWNGATE(ALTBAR^.COMBINATION, (CONNECTIVITY<>NOTSET));
3182             PARALTLIST := NEXT
3183         END;
3184     POSITIONSET(CHILDGATEPOSITIONS, LIMITPOSITIONS, CONNECTIVITY, 6);
3185     WHILE GATEBOTTOMS <> NIL DO
3186         BEGIN
3187             WITH GATEBOTTOMS^.GBOTTOM^ DO
3188                 BEGIN
3189                     CHILDPOSITIONS := CHILDGATEPOSITIONS;
3190                     IF CHILDPOSITIONS.COMBINED AND (CONNBONDS.CONNECTIONS = 1)
3191                         THEN ALTERCONNBONDS(CONNBONDS);
3192                     PG := CHILDPS^.PARENTGATE
3193                 END;
3194             WHILE PG <> NIL DO WITH PG^ DO
3195                 BEGIN
3196                     CHILDPOSITIONS := CHILDGATEPOSITIONS;

```

```

3197         IF CHILDPOSITIONS.COMBINED AND (CONNBONDS.CONNECTIONS=1)
3198             THEN ALTERCONNBONDS(CONNBONDS);
3199         PG := NEXT
3200     END;
3201     PTR := GATEBOTTOMS^.NEXT;
3202     DISPOSE(GATEBOTTOMS);
3203     GATEBOTTOMS := PTR
3204 END;
3205 NEXTTOKEN
3206 END;          { of MODIFYCHILDPOSITIONS
-----}
3207
3208
3209
3210
3211 {-----}
3212             PROCEDURE ELEMENT }
3213
3214 PROCEDURE ELEMENT(PARALTLIST      : PPALTBARS;
3215                   OPTIONALSUB    : BOOLEAN);
3216
3217 { Analyses a substituent definition element, building up the ECTR.
3218   Called by ALTNTVE}
3219
3220 VAR OPENERS,                {Valid tokens to begin an element}
3221     TERMINATORS             : DELIMSET;    {Valid tokens to end an element}
3222     VALID                   : BOOLEAN;
3223     DELIMCHECK              : DELIMTYPE;
3224     LIMITPOSITIONS         : TGROUPMEMS;
3225     GATEPARENTPOSITIONS    : PTGROUPMEMS; {Position set for inclusion in the gate}
3226     GATEFREQUENCY          : INTRECORD;   {Frequencies for inclusion in the gate}
3227     GATEPS                  : PTRPSTYPE;  {The child partial structure}
3228
3229
3230
3231 PROCEDURE SETPARENTGATE(COMBIN : PCOMBINLIST;
3232                        PARALT : PPALTBARS);
3233
3234 { Sets up a single new parent gate at the head of the list on COMBIN^.CHILDPS^.
3235   Called by SETCOMBARS
3236     SUBSTASVALUE}

```

```

3237
3238
3239     VAR NEWPG : PPARENTLIST;
3240
3241     BEGIN
3242     NEW(NEWPG);
3243     ECTRSIZE := ECTRSIZE + 26;
3244     NEWPG^.CHILDPOSITIONS := COMBIN^.CHILDPOSITIONS;
3245     IF GATEPARENTPOSITIONS=NIL
3246     THEN NEWPG^.PARENTPOSITIONS := PARALT^.PRNTPOSNS^
3247     ELSE NEWPG^.PARENTPOSITIONS := GATEPARENTPOSITIONS^;
3248     NEWPG^.PARENTPS := PARALT^.PARSTRUCT;
3249     NEWPG^.CONNBONDS := COMBIN^.CONNBONDS;
3250     NEWPG^.NEXT := COMBIN^.CHILDPS^.PARENTGATE;
3251     COMBIN^.CHILDPS^.PARENTGATE := NEWPG
3252     END;
3253
3254
3255
3256     FUNCTION NEWCOMBAR(PARALT      : PPALTBARS;
3257                       BARBOTTOM  : BOOLEAN) : PCOMBINLIST;
3258
3259     { Creates a combination bar, variant BARBOTTOM, and sets the non-variant fields,
3260     returning the bar as the result of the function.
3261     Called by SETCOMBARS
3262     EXTRALAYER
3263     NEWPARENTPSLIST}
3264
3265     VAR NEWCB : PCOMBINLIST;
3266
3267     BEGIN
3268     IF BARBOTTOM THEN NEW(NEWCB, TRUE)
3269     ELSE NEW(NEWCB, FALSE);
3270     ECTRSIZE := ECTRSIZE + 11 + (ORD(BARBOTTOM) * 13);
3271     WITH NEWCB^ DO
3272     BEGIN
3273     PARENTPOSITIONS := GATEPARENTPOSITIONS;
3274     FREQUENCY := GATEFREQUENCY;
3275     IF PARALT^.ALTBAR = NIL
3276     THEN BEGIN

```

```

3277      {There is no alternative bar. The combination list needs to be
3278      attached directly to the PS record, and PARENTPOSITIONS
3279      needs to be taken from PARALT^.PRNTPOSNS if none has been obtained
3280      from GATEPARENTPOSITIONS}
3281      NEXT := PARALT^.PARSTRUCT^.CHILDGATE;
3282      PARALT^.PARSTRUCT^.CHILDGATE := NEWCB;
3283      IF PARENTPOSITIONS = NIL THEN PARENTPOSITIONS := PARALT^.PRNTPOSNS
3284      END
3285      ELSE BEGIN
3286      {The combination list needs to be attached to the alternative bar}
3287      NEXT := PARALT^.ALTBAR^.COMBINATION;
3288      PARALT^.ALTBAR^.COMBINATION := NEWCB
3289      END;
3290      BOTTOMBAR := BARBOTTOM
3291      END;
3292      NEWCOMBAR := NEWCB
3293      END;
3294
3295
3296
3297      {-----}
3298      PROCEDURE SETCOMBARS(PARALTLIST : PPALTBARS;
3299                          GATEPS      : PTRPSTYPE);
3300
3301      { Sets up the child and parent gates for all the items in PARALTLIST
3302      Called by Body of ELEMENT}
3303
3304      VAR NEWCOMBIN : PCOMBINLIST;
3305
3306
3307
3308      FUNCTION NEEDTOCHECK(PARENTPS      : PTRPSTYPE;
3309                          NEWCONNBNDS,
3310                          OLDCONNBNDS : TCONNBNDS) : BOOLEAN;
3311
3312      BEGIN
3313      IF PARENTPS = NIL
3314      THEN NEEDTOCHECK := FALSE
3315      ELSE IF PARENTPS^.PSVARIETY = SPECIFIC
3316      THEN CASE NEWCONNBNDS.CONNECTIONS OF

```

```

3317         NOTSET,
3318         0       : PROGERROR(13);
3319         1       : CASE OLDCONNBONDS.CONNECTIONS OF
3320             NOTSET : NEEDTOCHECK := MAGNITUDE(NEWCONNBONDS.BOND) > 1;
3321             0, 2   : PROGERROR(14);
3322             1     : NEEDTOCHECK := MAGNITUDE(NEWCONNBONDS.BOND) > MAGNITUDE(OLDCONNBONDS.B
ND)

3323             END;
3324         2       : CASE OLDCONNBONDS.CONNECTIONS OF
3325             NOTSET : NEEDTOCHECK := (MAGNITUDE(NEWCONNBONDS.BONDA) > 1)
3326                 OR (MAGNITUDE(NEWCONNBONDS.BONDB) > 1);
3327             0, 1   : PROGERROR(15);
3328             2     : NEEDTOCHECK := (MAGNITUDE(NEWCONNBONDS.BONDA) > MAGNITUDE(OLDCONNBONDS
BONDA))
OR (MAGNITUDE(NEWCONNBONDS.BONDB) > MAGNITUDE(OLDCONNBONDS.BONDB))

3329
3330
3331             END
3332         END
3333     ELSE NEEDTOCHECK := FALSE
3334 END;

3335
3336
3337
3338     FUNCTION ORIGINALPOSNS(PARENTPS      : PTRPSTYPE;
3339                            LASTPARPOSNS : PTGROUPMEMS) : BOOLEAN;
3340
3341     { Returns TRUE if LASTPARPOSNS is the PRNTPOSNS field of any of the child gates
3342       leading from PARENTPS}
3343
3344     VAR CGPTR : PCOMBINLIST;
3345         ROWNO : ATOMNUMBER;
3346         FOUND : BOOLEAN;
3347
3348     BEGIN
3349         FOUND := FALSE;
3350         CGPTR := PARENTPS^.CHILDGATE;
3351         WHILE (CGPTR <> NIL) AND NOT FOUND DO WITH CGPTR^ DO
3352             BEGIN
3353                 FOUND := (PARENTPOSITIONS = LASTPARPOSNS);
3354                 CGPTR := NEXT

```

APPENDIX 3:

GENSAL INTERPRETER

```

3355     END;
3356     IF NOT FOUND THEN WITH PARENTPS^ DO
3357         FOR ROWNO := 1 TO MAXCT DO IF CT[ROWNO] <> NIL
3358             THEN WITH CT[ROWNO]^ DO IF NOT ATOMICROW
3359                 THEN BEGIN
3360                     CGPTR := VALUES;
3361                     WHILE (CGPTR <> NIL) AND NOT FOUND DO WITH CGPTR^ DO
3362                         BEGIN
3363                             FOUND := (PARENTPOSITIONS = LASTPARPOSNS);
3364                             CGPTR := NEXT
3365                         END
3366                     END;
3367     ORIGINALPOSNS := FOUND
3368     END;
3369
3370
3371
3372     PROCEDURE CHECKPOSNS(VAR BOTPARPOSNS : PTGROUPMEMS;
3373                         LASTPARPOSNS   : PTGROUPMEMS;
3374                         PARENTPS       : PTRPSTYPE;
3375                         BOND MAG       : TBOND MAG);
3376
3377     { Checks that any positions given in position sets are actually available for
3378       a bond of the MAGNITUDE in question, as this was not previously known. }
3379
3380     VAR NOPOSNSGIVEN : BOOLEAN;
3381         GIVENPOSNS   : INTEGSET;
3382         POSN         : ATOMNUMBER;
3383
3384     BEGIN
3385     IF BOTPARPOSNS = NIL
3386     THEN BEGIN
3387         NOPOSNSGIVEN := ORIGINALPOSNS(PARENTPS, LASTPARPOSNS);
3388         IF NOPOSNSGIVEN
3389         THEN BEGIN
3390             NEW(BOTPARPOSNS);
3391             ECTRSIZE := ECTRSIZE + 9;
3392             BOTPARPOSNS^ := LASTPARPOSNS^
3393         END;
3394     GIVENPOSNS := LASTPARPOSNS^.MEMBERS

```

```

3395         END
3396     ELSE BEGIN
3397         NOPOSNSGIVEN := FALSE;
3398         GIVENPOSNS := BOTPARPOSNS^.MEMBERS
3399     END;
3400     FOR POSN := 1 TO MAXCT DO IF POSN IN GIVENPOSNS
3401     THEN WITH PARENTPS^.CT[POSN]^ DO IF ATOMICROW
3402     THEN IF HYDROGENS < BOND MAG
3403     THEN IF NOPOSNSGIVEN
3404     THEN WITH BOTPARPOSNS^ DO MEMBERS := MEMBERS - [POSN]
3405     ELSE FAILURE(45, POSN, ' ');
3406     IF BOTPARPOSNS^.MEMBERS = [] THEN FAILURE(46, 0, ' ');
3407     END;
3408
3409
3410
3411     FUNCTION COMBINEDPOSITIONS(CONNBONDS : TCONNBONDS;
3412     LIMITPOSNS : INTEGSET;
3413     PARENTPS : PTRPSTYPE) : PTGROUPMEMS;
3414
3415     { Returns a COMBINED position set based on the available positions in PARENTPS.
3416     If the bond magnitudes are 1, then the previously-determined LIMITPOSNS can
3417     be used instead of calling GETAVAILABLEPOSITIONS.}
3418
3419     VAR COMBPOSNS : PTGROUPMEMS;
3420     MAGNIT : INTEGER;
3421     POSNSETA,
3422     POSNSETB,
3423     POSNSETC : INTEGSET;
3424
3425     BEGIN
3426     NEW(COMBPOSNS);
3427     ECTRSIZE := ECTRSIZE + 9;
3428     COMBPOSNS^.COMBINED := TRUE;
3429     COMBPOSNS^.COMBMEMS := NIL;
3430     MAGNIT := MAGNITUDE(CONNBONDS.BONDA);
3431     IF MAGNIT > 1 THEN GETAVAILABLEPOSITIONS(PARENTPS, POSNSETA, MAGNIT)
3432     ELSE POSNSETA := LIMITPOSNS;
3433     MAGNIT := MAGNITUDE(CONNBONDS.BONDB);
3434     IF MAGNIT > 1 THEN GETAVAILABLEPOSITIONS(PARENTPS, POSNSETB, MAGNIT)

```

```

3435             ELSE POSNSETB := LIMITPOSNS;
3436 MAGNIT := MAGNIT + MAGNITUDE(CONNBONDS.BONDA);
3437 IF MAGNIT <= 3 THEN GETAVAILABLEPOSITIONS(PARENTPS, POSNSETC, MAGNIT)
3438             ELSE POSNSETC := [];
3439 POSNSETC := POSNSETA * POSNSETB * POSNSETC;
3440 LISTPOSNS(COMBPOSNS^.COMBMEMS, POSNSETA, POSNSETB, POSNSETC);
3441 IF COMBPOSNS^.COMBMEMS = NIL THEN FAILURE(47, 0, ' ');
3442 COMBINEDPOSITIONS := COMBPOSNS
3443 END;
3444
3445
3446
3447 PROCEDURE CHECKCOMBPOSNS(VAR BOTPARPOSNS : PTGROUPMEMS;
3448                         LASTPARPOSNS   : PTGROUPMEMS;
3449                         PARENTPS       : PTRPSTYPE;
3450                         MAGA,
3451                         MAGB           : TBOND MAG);
3452
3453 { Checks that any positions specified in position sets are actually available
3454   for bonds of the MAGNITUDEs in question, which were not previously known.}
3455
3456 VAR REMOVEA,
3457     REMOVEB      : INTEGSET;
3458     NOPOSNSGIVEN : BOOLEAN;
3459     GIVENPOSNS,
3460     LISTPTR,
3461     DELPTR       : PDOUBLIST;
3462
3463 BEGIN
3464 IF BOTPARPOSNS= NIL
3465     THEN BEGIN
3466         NOPOSNSGIVEN := ORIGINALPOSNS(PARENTPS, LASTPARPOSNS);
3467         GIVENPOSNS := COPYLIST(LASTPARPOSNS^.COMBMEMS)
3468     END
3469     ELSE BEGIN
3470         NOPOSNSGIVEN := FALSE;
3471         GIVENPOSNS := COPYLIST(BOTPARPOSNS^.COMBMEMS)
3472     END;
3473 REMOVEA := [];
3474 REMOVEB := [];

```

```

3475 LISTPTR := GIVENPOSNS;
3476 WHILE LISTPTR <> NIL DO
3477     BEGIN
3478         WITH PARENTPS^.CT[LISTPTR^.FIRST]^ DO IF ATOMICROW
3479             THEN IF HYDROGENS < MAGA
3480                 THEN IF NOPOSNSGIVEN
3481                     THEN REMOVEA := REMOVEA + [LISTPTR^.FIRST]
3482                     ELSE FAILURE(45, LISTPTR^.FIRST, ' ');
3483         WITH PARENTPS^.CT[LISTPTR^.SECOND]^ DO IF ATOMICROW
3484             THEN IF HYDROGENS < MAGB
3485                 THEN IF NOPOSNSGIVEN
3486                     THEN REMOVEB := REMOVEB + [LISTPTR^.SECOND]
3487                     ELSE FAILURE(45, LISTPTR^.SECOND, ' ');
3488         LISTPTR := LISTPTR^.NEXT
3489     END;
3490 IF REMOVEA + REMOVEB = []
3491     THEN BEGIN
3492         REDUCEECTR(GIVENPOSNS);
3493         DESTROY(GIVENPOSNS)
3494     END
3495 ELSE BEGIN
3496     LISTPTR := GIVENPOSNS;
3497     WHILE LISTPTR <> NIL DO IF (LISTPTR^.FIRST IN REMOVEA) OR (LISTPTR^.SECOND IN REMOVEB)
3498         THEN BEGIN
3499             IF LISTPTR = GIVENPOSNS
3500                 THEN GIVENPOSNS := LISTPTR^.NEXT;
3501             DELPTR := LISTPTR;
3502             LISTPTR := LISTPTR^.NEXT;
3503             DISPOSE(DELPTR);
3504             ECTRSIZE := ECTRSIZE - 6
3505         END
3506     ELSE LISTPTR := LISTPTR^.NEXT;
3507 IF GIVENPOSNS = NIL THEN FAILURE(48, 0, ' ');
3508 IF BOTPARPOSNS = NIL
3509     THEN BEGIN
3510         NEW(BOTPARPOSNS);
3511         ECTRSIZE := ECTRSIZE + 9;
3512         BOTPARPOSNS^.COMBINED := TRUE
3513     END;
3514 BOTPARPOSNS^.COMBMEMS := GIVENPOSNS

```

```

3515         END
3516     END;
3517
3518
3519
3520     BEGIN {Body of SETCOMBARS}
3521     ADDTOLIST(GATEPS);
3522     WHILE PARALTLIST <> NIL DO
3523     BEGIN
3524         NEWCOMBIN := NEWCOMBAR(PARALTLIST, FALSE);
3525         IF PARALTLIST^.COPYCHILDPS
3526             THEN NEWCOMBIN^.CHILDPS := COPYPS(GATEPS)
3527             ELSE NEWCOMBIN^.CHILDPS := GATEPS;
3528         NEWCOMBIN^.CONNBONDS := PARALTLIST^.CONNBONDS;
3529         GETCHILDPOSITIONS(GATEPS,NEWCOMBIN^. CONNBONDS, NEWCOMBIN^.CHILDPOSITIONS);
3530     CASE NEWCOMBIN^.CONNBONDS.CONNECTIONS OF
3531     NOTSET : PROGERR(16);
3532     0      : ;
3533     1      : IF NEEDTOCHECK(PARALTLIST^.PARSTRUCT, NEWCOMBIN^.CONNBONDS, PARALTLIST^.CONNBONDS)
3534             THEN CHECKPOSNS(NEWCOMBIN^.PARENTPOSITIONS,
3535                             PARALTLIST^.PRNTPOSNS,
3536                             PARALTLIST^.PARSTRUCT,
3537                             MAGNITUDE(NEWCOMBIN^.CONNBONDS.BOND));
3538     2      : IF (NEWCOMBIN^.PARENTPOSITIONS = NIL) AND NOT PARALTLIST^.PRNTPOSNS^.COMBINED
3539             THEN {no position set can have been given, as it would have had to have been
3540                 COMBINED. Consequently we can reestablish PARENTPOSITIONS from scratch}
3541                 IF PARALTLIST^.PARSTRUCT <> NIL THEN
3542                     NEWCOMBIN^.PARENTPOSITIONS := COMBINEDPOSITIONS(NEWCOMBIN^.CONNBONDS,
3543                             PARALTLIST^.PRNTPOSNS^.MEMBERS,
3544                             PARALTLIST^.PARSTRUCT)
3545                 ELSE {any position sets we have must be COMBINED}
3546                     IF NEEDTOCHECK(PARALTLIST^.PARSTRUCT, NEWCOMBIN^.CONNBONDS, PARALTLIST^.CONNBONDS)
3547                         THEN CHECKCOMBPOSNS(NEWCOMBIN^.PARENTPOSITIONS,
3548                             PARALTLIST^.PRNTPOSNS,
3549                             PARALTLIST^.PARSTRUCT,
3550                             MAGNITUDE(NEWCOMBIN^.CONNBONDS.BONDA),
3551                             MAGNITUDE(NEWCOMBIN^.CONNBONDS.BONDB));
3552     END;
3553     IF NOT DEFNTABLEENTRY(PARALTLIST^.PARSTRUCT)
3554     THEN SETPARENTGATE(NEWCOMBIN, PARALTLIST);

```

APPENDIX 3:

GENSYM  
LINTERPRETER

```

3555     PARALTLIST := PARALTLIST^.NEXT
3556     END
3557     END;      {of SETCOMBARS
3558     -----}
3559
3560
3561     FUNCTION NEWPARENTPSLIST(PARALTLIST : PPALTBARS): PPSLIST;
3562
3563     { Sets up a new PPSLIST based on the items in PARALTLIST, which can be passed
3564       in a recursive call to ALTNVLIST. This procedure also establishes non-BOTTOMBAR
3565       combination bars in the gates for all items in PARALTLIST.
3566       Called from TRANSLATENOMEN
3567         Body of ELEMENT}
3568
3569     VAR LISTPTR,           { Points to the top of the growing list }
3570         WRITEPTR : PPSLIST; { New addition to the list }
3571         NEWCOMBIN : PCOMBINLIST; { New combination gate }
3572
3573     BEGIN
3574     LISTPTR := NIL;
3575     WHILE PARALTLIST <> NIL DO
3576     BEGIN
3577         NEWCOMBIN := NEWCOMBAR(PARALTLIST, TRUE);
3578         NEWCOMBIN^.ALTERNATIVES := NIL;
3579         NEW(WRITEPTR);
3580         WRITEPTR^.NEXT := LISTPTR;
3581         WRITEPTR^.PARSTRUCT := PARALTLIST^.PARSTRUCT;
3582         WRITEPTR^.CONNBONDS := PARALTLIST^.CONNBONDS;
3583         IF GATEPARENTPOSITIONS = NIL
3584         THEN WRITEPTR^.PRNTPOSNS := PARALTLIST^.PRNTPOSNS
3585         ELSE WRITEPTR^.PRNTPOSNS := GATEPARENTPOSITIONS;
3586         WRITEPTR^.COMBINS := NEWCOMBIN;
3587         WRITEPTR^.COPYCHILDPS := PARALTLIST^.COPYCHILDPS;
3588         WRITEPTR^.FURTHERSUB := NIL;
3589         LISTPTR := WRITEPTR;
3590         PARALTLIST := PARALTLIST^.NEXT
3591     END;
3592     NEWPARENTPSLIST := LISTPTR
3593     END;
3594

```

```

3595
3596
3597 PROCEDURE GETSPSPARAMS(VAR GATEPS : PTRPSTYPE);
3598
3599 { Sets up a GENERIC PS, and initialises the PARAMLIST with the global
3600 SPSPARAMLIST, which has been set up by SPSVARIETY.
3601 Called by TRANSLATENOMEN}
3602
3603 VAR PARAM : TPARAMETERS;
3604 PTR : PDOUBLIST;
3605
3606 BEGIN
3607 NEW(GATEPS, GENERIC);
3608 ECTRSIZE := ECTRSIZE + 50;
3609 WITH GATEPS^ DO
3610 BEGIN
3611 PSVARIETY := GENERIC;
3612 VISITED := FALSE;
3613 PARENTGATE := NIL;
3614 CHILDGATE := NIL;
3615 PARAMLIST := SPSPARAMLIST;
3616 FOR PARAM := ATOMCOUNT TO HETEROATOM DO
3617 BEGIN
3618 PTR := PARAMLIST[PARAM].SUBRANGES;
3619 WHILE PTR <> NIL DO
3620 BEGIN
3621 ECTRSIZE := ECTRSIZE + 6;
3622 PTR := PTR^.NEXT
3623 END
3624 END
3625 END
3626 END;
3627
3628
3629
3630 {-----}
3631 PROCEDURE PARAMETERLIST(GATEPS : PTRPSTYPE);
3632
3633 { Analyses a Gensal parameter list. For standard parameters, the information
3634 is stored in the appropriate element of PARAMLIST, this being determined

```

```

3635     by the function PARAMETER. The existing parameter values are used
3636     to limit those analysed, by creating LIMITSET, to be passed to SELECTOR.
3637     In SELECTOR this is converted back to an INTRECORD for passing to
3638     INTEGERRANGE, but this is not as silly as it seems, as otherwise
3639     INTEGERRANGE would be using the same INTRECORD linked list both as
3640     LIMITRANGE and RANGEVALUES. If GATEPS is not GENERIC then the values for
3641     each parameter are placed in DUMMYPARAM, and any PDOUBLIST linked list
3642     immediately DESTROYed. Non-standard parameters are handled by USERPARAMETER.
3643     Called by TRANSLATENOMEN}
3644
3645     VAR PARAMIDS    : DELIMSET;
3646         DELIMCHECK : DELIMTYPE;
3647         DUMMYPARAM  : INTRECORD;
3648         LIMITSET    : INTEGSET;
3649
3650
3651
3652     FUNCTION PARAMETER(PARAMDELIM : DELIMTYPE) : TPARAMETERS;
3653
3654     { Returns the PARAMETER that is equivalent to the delimiter passed as PARAMDELIM}
3655
3656     BEGIN
3657     CASE PARAMDELIM OF
3658         GC : PARAMETER := ATOMCOUNT;
3659         GT : PARAMETER := TBRANCH;
3660         GQ : PARAMETER := QBRANCH;
3661         GE : PARAMETER := EUNSATURATION;
3662         GY : PARAMETER := YUNSATURATION;
3663         GRC : PARAMETER := RINGCOUNT;
3664         GRN : PARAMETER := RINGATOMS;
3665         GRS : PARAMETER := RINGSUBSTITUTION;
3666         GRF : PARAMETER := RINGFUSIONS;
3667         GRA : PARAMETER := RINGAROMATIC;
3668         GZ : PARAMETER := HETEROATOM
3669     END;
3670 END;
3671
3672
3673
3674     PROCEDURE FINDCONNECTIONS(VAR CONNBONDS : TCONNBONDS;

```

```

3675         VAR PRNTPOSNS : PTGROUPMEMS;
3676         PSADDRESS      : PTRPSTYPE;
3677         SUBST           : SUBSTITUENT);
3678
3679 { Sets CONNBONDS and PRNTPOSNS for the declaration of SUBST referencing
3680 PSADDRESS. Any previous declaration of SUBST will give a value for
3681 CONNBONDS - otherwise CONNBONDS.CONNECTIONS is set to NOTSET.
3682 Since PSADDRESS^.PSVARIETY is GENERIC or OTHER, the bond magnitude is
3683 irrelevant to the determination of the available positions - hence a
3684 dummy value of 1 is passed to GETAVAILABLEPOSITIONS.
3685 Called by USERPARAMETER}
3686
3687     VAR AVAILPOSNS : INTEGSET;
3688
3689     BEGIN
3690     IF SUBST IN DECLSUBS THEN CONNBONDS := RDECLARATIONTABLE[SUBST]^ .CONNBONDS
3691     ELSE CONNBONDS.CONNECTIONS := NOTSET;
3692
3693     NEW(PRNTPOSNS);
3694     ECTRSIZE := ECTRSIZE + 9;
3695     PRNTPOSNS^.COMBINED := CONNBONDS.CONNECTIONS = 2;
3696     CASE CONNBONDS.CONNECTIONS OF
3697     0       : PRNTPOSNS^.MEMBERS := [];
3698     NOTSET,
3699     1       : GETAVAILABLEPOSITIONS(PSADDRESS, PRNTPOSNS^.MEMBERS, 1);
3700     2       : BEGIN
3701               GETAVAILABLEPOSITIONS(PSADDRESS, AVAILPOSNS, 1);
3702               PRNTPOSNS^.COMBMEMS := NIL;
3703               LISTPOSNS(PRNTPOSNS^.COMBMEMS, AVAILPOSNS+[0], AVAILPOSNS+[0], AVAILPOSNS)
3704     END
3705     END
3706     END;
3707
3708
3709     PROCEDURE USERPARAMETER(GATEPS : PTRPSTYPE);
3710
3711     { Analyses a user-defined parameter. FINDCONNECTIONS is used to determine
3712     CONNBONDS and PARENTPOSITIONS, for passing to DECLARESUBST.}
3713
3714     VAR NEWCOMBIN      : PCOMBINLIST;

```

APPENDIX 3:

GENSAL INTERPRETER

```

3715         CONNBONDS      : TCONNBONDS;
3716         PARENTPOSITIONS : PTGROUPMEMS;
3717
3718     BEGIN
3719     NEXTTOKEN;
3720     WHILE CHECKDELIM([GR])=INVALIDTOKEN DO ERROR(26,0);
3721     NEXTTOKEN;
3722     REPEAT
3723     WHILE TOKEN.NATURE<>INTEGRAL DO ERROR(23,0);
3724     IF NOT (TOKEN.INTEGVAL IN [1..MAXVARS]) THEN ERROR(28,0)
3725     UNTIL TOKEN.INTEGVAL IN [1..MAXVARS];
3726     FINDCONNECTIONS(CONNBONDS, PARENTPOSITIONS, GATEPS, TOKEN.INTEGVAL);
3727     DECLARESUBST(TOKEN.INTEGVAL,      {subst}
3728                 GATEPS,              {psaddress}
3729                 NIL,                  {savps}
3730                 CONNBONDS,            {connbonds}
3731                 PARENTPOSITIONS);    {prntposns}
3732     ENTERCOMBIN(TOKEN.INTEGVAL, GATEPS^.CHILDGATE);
3733     NEXTTOKEN;
3734     WHILE CHECKDELIM([GPRIME]) = INVALIDTOKEN DO ERROR(29,0);
3735     NEXTTOKEN;
3736     SELECTOR(GATEPS^.CHILDGATE^.FREQUENCY, [0..MAXVARS], 5)
3737     END;
3738
3739
3740
3741     BEGIN {Body of PARAMETERLIST}
3742     PARAMIDS := [GC, GT, GQ, GE, GY, GRC, GRN, GRS, GRF, GRA, GZ];
3743     REPEAT
3744     DELIMCHECK := CHECKDELIM([GOPENANG, GLSQUARE, GPRIME]+PARAMIDS+TERMINATORS);
3745     IF DELIMCHECK=INVALIDTOKEN THEN ERROR(24,0)
3746     UNTIL DELIMCHECK <> INVALIDTOKEN;
3747     IF DELIMCHECK = GOPENANG
3748     THEN DELIMCHECK := GC {Parameter identifier not needed}
3749     ELSE IF DELIMCHECK IN PARAMIDS THEN NEXTTOKEN;
3750     WHILE DELIMCHECK IN PARAMIDS+[GPRIME] DO
3751     BEGIN
3752     IF DELIMCHECK = GPRIME
3753     THEN USERPARAMETER(GATEPS)
3754     ELSE IF GATEPS^.PSVARIETY = GENERIC

```

```

3755         THEN BEGIN
3756             INTSET(LIMITSET, GATEPS^.PARAMLIST[PARAMETER(DELIMCHECK)]);
3757             SELECTOR(GATEPS^.PARAMLIST[PARAMETER(DELIMCHECK)], LIMITSET, 5)
3758         END
3759         ELSE BEGIN
3760             SELECTOR(DUMMYPARAM, [0..MAXVARS], 5);
3761             REDUCECTR(DUMMYPARAM.SUBRANGES);
3762             DESTROY(DUMMYPARAM.SUBRANGES)
3763         END;
3764     NEXTTOKEN;
3765     REPEAT
3766         DELIMCHECK := CHECKDELIM([GPRIME, GLSQUARE]+PARAMIDS+TERMINATORS);
3767         IF DELIMCHECK=INVALIDTOKEN THEN ERROR(24,0)
3768         UNTIL DELIMCHECK <> INVALIDTOKEN;
3769         IF DELIMCHECK IN PARAMIDS THEN NEXTTOKEN
3770     END
3771 END;
3772     { of PARAMETERLIST
3773 -----}
3774
3775
3776
3777 {-----}
3778 PROCEDURE TRANSLATENOMEN(VAR GATEPS : PTRPSTYPE);
3779
3780 { Determines whether or not a record is held for the current TOKEN.NOMENVAL
3781   (taking synonyms in RECORDHELD, if terminal input), and sets up an
3782   appropriate child PS, or pushes down the input environment and makes a
3783   recursive call to ALTNVLIST. If no record is held, then an OTHER PS is set
3784   up.
3785   Called from body of ELEMENT}
3786
3787 VAR ADDRESS          : INTEGER;          { SPSfile address for TOKEN.NOMENVAL }
3788     OLDN              : 0..MAXLENGTH;    { Saved value of N }
3789     OLDBUFFER         : LINESTRING;      { Saved value of BUFFER }
3790     INSERTLINES,     :                   { Lines being inserted }
3791     OLDCURRENTLINE   : PLINELIST;       { Saved value of CURRENTLINE }
3792     OLDMODE           : TINPUTMODE;     { Saved INPUTMODE }
3793
3794

```

APPENDIX 3:

GENSAL INTERPRETER

```

3795
3796 PROCEDURE MODIFYGATEPOSITIONS(COMBINBAR      : PCOMBINLIST;
3797                               HSTAVAILPOSNS : PTGROUPMEMS);
3798
3799 { Traces down a child gate (headed by COMBINBAR), altering the PARENTPOSITIONS
3800   field in each, to conform to HSTAVAILPOSNS. If the bar is BOTTOMBAR
3801   then the PARENTPOSITIONS fields of the corresponding parent gates from the
3802   CHILDPS are altered similarly; otherwise the procedure recurses on itself
3803   for each COMBINATION in the ALTERNATIVES.}
3804
3805 VAR ALTERNBAR : PALTERNLIST;
3806     PARENTGATE : PPARENTLIST;
3807     TOPBAR     : BOOLEAN;
3808     POSN      : ATOMNUMBER;
3809
3810 BEGIN
3811 TOPBAR := (HSTAVAILPOSNS = NIL);
3812 IF TOPBAR THEN
3813   BEGIN
3814     NEW(HSTAVAILPOSNS);
3815     WITH HSTAVAILPOSNS^ DO
3816       BEGIN
3817         COMBINED := FALSE;
3818         GETAVAILABLEPOSITIONS(INSERTHSTPS, MEMBERS, 1);
3819         MEMBERS := MEMBERS + [0]
3820       END
3821     END;
3822 WHILE COMBINBAR <> NIL DO WITH COMBINBAR^ DO
3823   BEGIN
3824     IF PARENTPOSITIONS <> NIL THEN WITH PARENTPOSITIONS^ DO
3825       IF TOPBAR
3826         THEN IF COMBINED
3827           THEN BEGIN
3828             REDUCE(COMBMEMS, HSTAVAILPOSNS^);
3829             IF COMBMEMS=NIL THEN FAILURE(4, 0, '  ')
3830           END
3831           ELSE BEGIN
3832             MEMBERS := MEMBERS * HSTAVAILPOSNS^.MEMBERS;
3833             IF MEMBERS = [] THEN FAILURE(4, 0, '  ')
3834           END

```

```

3835         ELSE IF COMBINED
3836             THEN CHECKALLWITHIN(COBBMEMS, HSTAVAILPOSNS^.MEMBERS, 3)
3837             ELSE IF MEMBERS <= HSTAVAILPOSNS^.MEMBERS
3838                 THEN {OK}
3839                 ELSE FOR POSN := 1 TO MAXCT DO
3840                     IF (POSN IN MEMBERS) AND NOT (POSN IN HSTAVAILPOSNS^.MEMBERS)
3841                         THEN FAILURE(3, POSN, ' ');
3842     IF BOTTOMBAR
3843     THEN BEGIN
3844         PARENTGATE := CHILDPS^.PARENTGATE;
3845         WHILE PARENTGATE <> NIL DO WITH PARENTGATE^ DO
3846             BEGIN
3847                 IF PARENTPS = INSERTHSTPS
3848                     THEN WITH PARENTPOSITIONS DO
3849                         IF COMBINED
3850                             THEN {no need to do anything - pointers point to same PDOUBLIST in parent a
3851                                 d child gates}
3852                             ELSE MEMBERS := MEMBERS * HSTAVAILPOSNS^.MEMBERS;
3853                             PARENTGATE := NEXT
3854                             END
3855             END
3856         ELSE BEGIN
3857             ALTERNBAR := ALTERNATIVES;
3858             WHILE ALTERNBAR <> NIL DO WITH ALTERNBAR^ DO
3859                 BEGIN
3860                     MODIFYGATEPOSITIONS(COMBINATION, HSTAVAILPOSNS);
3861                     ALTERNBAR := NEXT
3862                     END;
3863             COMBINBAR := NEXT
3864             END;
3865     IF TOPBAR THEN DISPOSE(HSTAVAILPOSNS)
3866
3867     END;
3868
3869
3870
3871     BEGIN { Body of Procedure TRANSLATENOMEN }
3872     IF RECORDHELD(TOKEN.NOMENVAL, ADDRESS)
3873         THEN CASE SPSVARIETY(ADDRESS, FALSE) OF

```

APPENDIX 3:

GENSAL INTERPRETER

```

3874     SPECIFIC : BEGIN
3875         NEW(GATEPS, SPECIFIC);
3876         ECTRSIZE := ECTRSIZE + 70;
3877         WITH GATEPS^ DO
3878             BEGIN
3879                 PSVARIETY := SPECIFIC;
3880                 VISITED := FALSE;
3881                 PARENTGATE := NIL;
3882                 CHILDGATE := NIL;
3883                 PROCESSCT(CT, FALSE, GATEPS)
3884             END;
3885         NEXTTOKEN
3886     END;
3887
3888     GENERIC  : BEGIN
3889         GETSPSPARAMS(GATEPS);
3890         NEXTTOKEN;
3891         IF CHECKDELIM([GLSQUARE]+TERMINATORS)=INVALIDTOKEN
3892             THEN PARAMETERLIST(GATEPS);
3893         IF INPUTMODE = INSERTTEXT
3894             THEN IF INSERTHSTPS = NIL
3895                 THEN INSERTHSTPS := GATEPS
3896                    ELSE PROGERROR(17) { multiple HSTs in SPSfile expression}
3897         END;
3898
3899     OTHER   : BEGIN
3900         GATEPS := NIL;
3901         { Save current environment }
3902         OLDCURRENTLINE := CURRENTLINE;
3903         CURRENTLINE := INSERTGENEX;
3904         OLDMODE := INPUTMODE;
3905         INPUTMODE := INSERTTEXT;
3906         OLDBUFFER := BUFFER;
3907         OLDN := N;
3908         N := MAXLENGTH;
3909         IF OLDMODE <> INSERTTEXT THEN INSERTHSTPS := NIL;
3910         INSERTLINES := INSERTGENEX;
3911         ALTNVLIST(NEWPARENTPSLIST(PARALTLIST), FALSE);
3912         IF CHECKDELIM( [GRPAREN])=INVALIDTOKEN THEN
3913             PROGERROR(18); {missing "" in SPSfile expression}

```

APPENDIX 3:

GENSAL INTERPRETER

```

3914      { Restore former environment }
3915      CURRENTLINE := OLDCURRENTLINE;
3916      BUFFER := OLDBUFFER;
3917      INPUTMODE := OLDMODE;
3918      N := OLDN;
3919      DELETEGENSAL(INSERTLINES);
3920      NEXTTOKEN;
3921      IF (INPUTMODE <> INSERTTEXT) AND (INSERTHSTPS <> NIL)
3922      THEN IF CHECKDELIM([GLSQUARE]+TERMINATORS)=INVALIDTOKEN
3923      THEN BEGIN
3924          PARAMETERLIST(INSERTHSTPS);
3925          IF INSERTHSTPS^.CHILDGATE <> NIL THEN
3926              MODIFYGATEPOSITIONS(INSERTHSTPS^.CHILDGATE, NIL)
3927      END
3928      END
3929      ELSE BEGIN
3930          NEW(GATEPS, OTHER);
3931          ECTRSIZE := ECTRSIZE + 22;
3932          WITH GATEPS^ DO
3933              BEGIN
3934                  PSVARIETY := OTHER;
3935                  VISITED := FALSE;
3936                  PARENTGATE := NIL;
3937                  CHILDGATE := NIL;
3938                  TERM := TOKEN.NOMENVAL
3939              END;
3940          NEXTTOKEN;
3941          IF CHECKDELIM([GLSQUARE]+TERMINATORS)=INVALIDTOKEN
3942          THEN PARAMETERLIST(GATEPS)
3943      END
3944      END;
3945      { of TRANSLATENOMEN
3946      -----}
3947
3948
3949
3950      FUNCTION EXTRALAYER(PARALT      : PPALTBARS;
3951                          DUMMYSAVPS : PTRPSTYPE) : PCOMBINLIST;
3952
3953      { Creates an extra layer in the gate of which PARALT gives an ALTERNATIVE bar,

```

```

3954         inserts one alternative into it, for DUMMYSAVPS.
3955         Called by SUBSTASVALUE}
3956
3957     VAR XLAYER : PCOMBINLIST;
3958
3959     BEGIN
3960     XLAYER := NEWCOMBAR(PARALT, TRUE);
3961     NEW(XLAYER^.ALTERNATIVES);
3962     ECTRSIZE := ECTRSIZE + 4;
3963     WITH XLAYER^, ALTERNATIVES^ DO
3964     BEGIN
3965     NEXT := NIL;
3966     NEW(COMBINATION, FALSE);
3967     ECTRSIZE := ECTRSIZE + 24;
3968     WITH COMBINATION^ DO
3969     BEGIN
3970     PARENTPOSITIONS := NIL;
3971     FREQUENCY.TOPRANGE := NOTSET;
3972     FREQUENCY.SUBRANGES := ESSENTFREQ;
3973     NEXT := NIL;
3974     BOTTOMBAR := TRUE;
3975     CHILDPS := DUMMYSAVPS;
3976     CONNBONDS := PARALT^.CONNBONDS;
3977     WITH CHILDPOSITIONS DO
3978     BEGIN
3979     COMBINED := CONNBONDS.CONNECTIONS = 2;
3980     IF COMBINED
3981     THEN BEGIN
3982     COMBMEMS := NIL;
3983     LISTPOSNS(COMBMEMS, [1..MAXCT], [1..MAXCT], [1..MAXCT])
3984     END
3985     ELSE MEMBERS := [1..MAXCT]
3986     END
3987     END
3988     END;
3989     EXTRALAYER := XLAYER
3990     END;
3991
3992
3993

```

```

3994 FUNCTION VALIDSUBST(PARALTLIST : PPALTBARS) : SUBSTITUENT;
3995
3996 { Obtains a substituent name and checks that it is in the range 1-63, and
3997   that if it is previously declared, its connectivity is compatible with
3998   that of the items in PARALTLIST.
3999   Called by SUBSTASVALUE}
4000
4001 LABEL 20;
4002
4003 BEGIN
4004 20 :
4005 WHILE TOKEN.NATURE <> INTEGRAL DO ERROR(23,0);
4006 IF NOT (TOKEN.INTEGVAL IN [1..MAXVARS])
4007   THEN BEGIN
4008     ERROR(28,0);
4009     GOTO 20
4010   END
4011 ELSE IF TOKEN.INTEGVAL IN DECLSUBS
4012   THEN WITH RDECLARATIONTABLE[TOKEN.INTEGVAL]^ DO
4013     IF CONNBONDS.CONNECTIONS = NOTSET
4014       THEN IF PARALTLIST^.CONNBONDS.CONNECTIONS = NOTSET
4015         THEN {no further information}
4016         ELSE BEGIN
4017           SETCONNBONDS(CONNBONDS, PARALTLIST^.CONNBONDS.CONNECTIONS);
4018           UPDATEPPSCONNS(RDECLARATIONTABLE[TOKEN.INTEGVAL])
4019         END
4020       ELSE IF PARALTLIST^.CONNBONDS.CONNECTIONS = NOTSET
4021         THEN BEGIN
4022           SETCONNBONDS(PARALTLIST^.CONNBONDS, CONNBONDS.CONNECTIONS);
4023           UPDATEPARALTCONNS(PARALTLIST)
4024         END
4025       ELSE IF CONNBONDS.CONNECTIONS = PARALTLIST^.CONNBONDS.CONNECTIONS
4026         THEN {compatible}
4027         ELSE BEGIN
4028           ERROR(31,0);
4029           GOTO 20
4030         END;
4031 VALIDSUBST := TOKEN.INTEGVAL
4032 END;
4033

```

```
4034
4035
4036 PROCEDURE SUBSTASVALUE(PARALTLIST : PPALTBARS);
4037
4038 { Analyses a substituent given as a substituent value, creates a DUMMY PS, and
4039   declares it using DECLARESUBST for each of the items in PARALTLIST. If
4040   the substituent has already been defined, copies the definition.
4041   Called from the body of ELEMENT}
4042
4043 VAR SUBST          : SUBSTITUENT;
4044     DUMMYCHILD,
4045     DUMMYSAVPS     : PTRPSTYPE;
4046     DECNPPPOSNS   : PTGROUPMEMS;
4047     PREVDEFN      : PALTERNLIST;
4048     OMITPG        : BOOLEAN;
4049     GATECONNBONDS : TCONNBONDS;
4050     DUMMYPOSNS    : TGROUPMEMS;
4051
4052 BEGIN
4053 NEXTTOKEN;
4054 SUBST := VALIDSUBST(PARALTLIST);
4055 NEW(DUMMYSAVPS, DUMMY);
4056 ECTRSIZE := ECTRSIZE + 8;
4057 WITH DUMMYSAVPS^ DO
4058 BEGIN
4059     PSVARIETY := DUMMY;
4060     VISITED := FALSE;
4061     CHILDGATE := NIL;
4062     PARENTGATE := NIL;
4063     SUBSTNAME := SUBST
4064 END;
4065 ADDTOLIST(DUMMYSAVPS);
4066 GATECONNBONDS := PARALTLIST^.CONNBONDS;
4067 GETCHILDPOSITIONS(DUMMYSAVPS, GATECONNBONDS, DUMMYPOSNS);
4068
4069 WHILE PARALTLIST <> NIL DO WITH PARALTLIST^ DO
4070 BEGIN
4071     IF GATEPARENTPOSITIONS=NIL
4072     THEN DECNPPPOSNS := PRNTPOSNS
4073     ELSE DECNPPPOSNS := GATEPARENTPOSITIONS;
```

```

4074     IF PARALTLIST^.COPYCHILDPS THEN DUMMYCHILD := COPYPS(DUMMYSAVPS)
4075                                     ELSE DUMMYCHILD := DUMMYSAVPS;
4076     DECLARESUBST(SUBST,
4077                 PARSTRUCT,
4078                 DUMMYCHILD,
4079                 CONNBONDS,
4080                 DECNPPOSNS);
4081     OMITPG := DEFNTABLEENTRY(PARSTRUCT);
4082     IF SUBST IN DEFNSUBS
4083     THEN PREVDEFN := RDEFINITIONTABLE[SUBST]^ALTERNATIVES
4084     ELSE PREVDEFN := NIL;
4085     RDECLARATIONTABLE[SUBST]^COMBINS := EXTRALAYER(PARALTLIST, DUMMYCHILD);
4086     IF NOT OMITPG THEN
4087     SETPARENTGATE(RDECLARATIONTABLE[SUBST]^COMBINS^.ALTERNATIVES^.COMBINATION, PARALTLIST);
4088     WHILE PREVDEFN <> NIL DO WITH RDECLARATIONTABLE[SUBST]^ DO
4089     BEGIN
4090         COPYALTBAR(COMBINS^.ALTERNATIVES, {newaltbar}
4091                 PREVDEFN^.COMBINATION, {oldcomblist}
4092                 COMBINS, {lastcomblayer}
4093                 PRNTPOSNS, {lastpposns}
4094                 PARSTRUCT, {prntps}
4095                 FALSE, {firstbar}
4096                 OMITPG, {omitpg}
4097                 TRUE); {copyyps}
4098         PREVDEFN := PREVDEFN^.NEXT
4099     END;
4100     PARALTLIST := NEXT
4101 END
4102 END;
4103
4104
4105
4106     FUNCTION NOVARIABLESUBTN(LIMITPOSITIONS : TGROUPMEMS) : BOOLEAN;
4107
4108     { True if LIMITPOSITIONS is empty.
4109     Called by Body of ELEMENT}
4110
4111     BEGIN
4112     WITH LIMITPOSITIONS DO
4113     IF COMBINED THEN NOVARIABLESUBTN := (COMBMEMS=NIL)

```

```

4114             ELSE NOVARIABLESUBTN := (MEMBERS=[])
4115 END;
4116
4117
4118
4119 PROCEDURE ADDZERO(VAR GATEFREQUENCY : INTRECORD);
4120
4121 { Adds zero to the integer range in GATEFREQUENCY.
4122   Called from body of ELEMENT}
4123
4124 VAR PTR : PDOUBLIST;
4125
4126 BEGIN
4127 WITH GATEFREQUENCY DO
4128   IF TOPRANGE = 0
4129     THEN { zero already present - no action needed }
4130     ELSE IF TOPRANGE = 1
4131       THEN TOPRANGE := 0
4132       ELSE IF SUBRANGES = NIL
4133         THEN SUBRANGES := ZEROFREQ
4134         ELSE BEGIN
4135           PTR := SUBRANGES;
4136           WHILE PTR^.NEXT <> NIL DO PTR := PTR^.NEXT;
4137           IF PTR^.FIRST = 0
4138             THEN { zero already present - no action needed }
4139             ELSE IF PTR^.FIRST = 1
4140               THEN PTR^.FIRST := 0
4141               ELSE PTR^.NEXT := ZEROFREQ
4142           END
4143   END;
4144
4145
4146
4147
4148
4149
4150 {-----}
4151 PROCEDURE GETLIMITPOSITIONS(PARALTLIST      : PPALTBARS;
4152                             VAR LIMITPOSITIONS : TGROUPMEMS);
4153

```

```

4154 { Establishes LIMITPOSITIONS from the contents of PARALTLIST.
4155   Called by Body of ELEMENT}
4156
4157 VAR ALLGENERIC : BOOLEAN;
4158
4159 BEGIN
4160 LIMITPOSITIONS.COMBINED := (PARALTLIST^.CONNBONDS.CONNECTIONS=2);
4161 WITH PARALTLIST^.PRNTPOSNS^ DO
4162   IF LIMITPOSITIONS.COMBINED
4163     THEN IF COMBINED
4164       THEN LIMITPOSITIONS.COMBMEMS := COPYLIST(COMBMEMS)
4165       ELSE BEGIN
4166         LIMITPOSITIONS.COMBMEMS := NIL;
4167         LISTPOSNS(LIMITPOSITIONS.COMBMEMS, MEMBERS, MEMBERS, MEMBERS)
4168       END
4169     ELSE IF COMBINED
4170       THEN PROGERROR(19) {combined position set with connectivity <> 2}
4171       ELSE LIMITPOSITIONS.MEMBERS := MEMBERS;
4172 WITH PARALTLIST^ DO
4173   BEGIN
4174     IF PARSTRUCT = NIL
4175       THEN ALLGENERIC := TRUE
4176       ELSE ALLGENERIC := (PARSTRUCT^.PSVARIETY = GENERIC);
4177     PARALTLIST := NEXT
4178   END;
4179 WHILE PARALTLIST <> NIL DO WITH PARALTLIST^ DO
4180   BEGIN
4181     IF ALLGENERIC AND (PARSTRUCT <> NIL)
4182       THEN ALLGENERIC := (PARSTRUCT^.PSVARIETY = GENERIC);
4183     IF LIMITPOSITIONS.COMBINED
4184       THEN REDUCE(LIMITPOSITIONS.COMBMEMS, PRNTPOSNS^)
4185       ELSE WITH PRNTPOSNS^ DO
4186         IF COMBINED
4187           THEN PROGERROR(20) { combined position set with connectivity <> 2}
4188           ELSE LIMITPOSITIONS.MEMBERS := LIMITPOSITIONS.MEMBERS * MEMBERS;
4189     PARALTLIST := NEXT
4190   END;
4191 WITH LIMITPOSITIONS DO
4192   IF ALLGENERIC AND NOT COMBINED
4193     THEN MEMBERS := MEMBERS + [0]

```

```

4194 END;
4195 {-----}
4196
4197
4198
4199 BEGIN {Body of Procedure ELEMENT }
4200 OPENERS := [GLSQUARE, GOPENANG, GR, GQUEST, GSD, GPRIME, GLPAREN];
4201 TERMINATORS := [GRPAREN, GAMPERSAND, GSB, GOSB, GSLASH, GSEMI, GELSE, GEND, GPERIOD, GAND, GOR];
4202 NEXTTOKEN;
4203 VALID := FALSE;
4204 REPEAT
4205     DELIMCHECK := CHECKDELIM(OPENERS);
4206     IF (DELIMCHECK=INVALIDTOKEN) AND (TOKEN.NATURE<>NOMENCLATURE)
4207         THEN ERROR(18,0)
4208         ELSE VALID := TRUE
4209 UNTIL VALID;
4210
4211 GATEPARENTPOSITIONS := NIL;
4212 IF DELIMCHECK = GLSQUARE THEN
4213     BEGIN
4214         GETLIMITPOSITIONS(PARALTLIST, LIMITPOSITIONS);
4215         IF NOVARIABLESUBTN(LIMITPOSITIONS)
4216             THEN WRITELN('Common definition of variable-substitution positions not possible.')
4217             ELSE BEGIN
4218                 NEW(GATEPARENTPOSITIONS);
4219                 ECTRSIZE := ECTRSIZE + 9;
4220                 POSITIONSET(GATEPARENTPOSITIONS^, LIMITPOSITIONS, PARALTLIST^.CONNBONDS.CONNECTIONS, 7);
4221                 IF PARALTLIST^.CONNBONDS.CONNECTIONS = NOTSET
4222                     THEN BEGIN
4223                         IF GATEPARENTPOSITIONS^.COMBINED
4224                             THEN SETCONNBONDS(PARALTLIST^.CONNBONDS, 2)
4225                             ELSE SETCONNBONDS(PARALTLIST^.CONNBONDS, 1);
4226                         UPDATEPARALTCONNS(PARALTLIST)
4227                     END;
4228                 NEXTTOKEN
4229             END;
4230         VALID := FALSE;
4231     REPEAT
4232         DELIMCHECK := CHECKDELIM(OPENERS-[GLSQUARE]);
4233         IF (DELIMCHECK=INVALIDTOKEN) AND (TOKEN.NATURE <> NOMENCLATURE)

```

APPENDIX

SENSAL INTERPRETER

```

4234         THEN ERROR(24,0)
4235         ELSE VALID := TRUE
4236     UNTIL VALID
4237     END;
4238
4239     IF DELIMCHECK=GOPENANG
4240     THEN BEGIN
4241         SELECTOR(GATEFREQUENCY, [0..MAXVARS], 8);
4242         IF OPTIONALSUB THEN ADDZERO(GATEFREQUENCY);
4243         NEXTTOKEN;
4244         VALID := FALSE;
4245         REPEAT
4246             DELIMCHECK := CHECKDELIM( OPENERS - [GLSQUARE, GOPENANG]);
4247             IF (DELIMCHECK=INVALIDTOKEN) AND (TOKEN.NATURE <> NOMENCLATURE)
4248                 THEN ERROR(24,0)
4249                 ELSE VALID := TRUE
4250             UNTIL VALID
4251         END
4252     ELSE WITH GATEFREQUENCY DO
4253         BEGIN
4254             TOPRANGE := NOTSET;
4255             IF OPTIONALSUB THEN SUBRANGES := OPTFREQ
4256                 ELSE SUBRANGES := ESSENTFREQ
4257         END;
4258
4259     CASE DELIMCHECK OF
4260
4261     GPRIME      : CONCATENATETERMS(GATEPS);
4262
4263     GQUEST     : BEGIN
4264                 NEW(GATEPS, UNKNOWN);
4265                 ECTRSIZE := ECTRSIZE + 6;
4266                 WITH GATEPS^ DO
4267                     BEGIN
4268                         PSVARIETY := UNKNOWN;
4269                         VISITED := FALSE;
4270                         PARENTGATE := NIL;
4271                         CHILDGATE := NIL
4272                     END
4273                 END;

```

```

4274
4275      GSD          : READSD(GATEPS, INPUTMODE=TERMINAL);
4276
4277      INVALIDTOKEN : TRANSLATENOMEN(GATEPS);
4278
4279      GR           : BEGIN
4280                  SUBSTASVALUE(PARALTLIST);
4281                  GATEPS := NIL
4282                  END;
4283
4284      GLPAREN      : BEGIN
4285                  ALTNVLIST(NEWPARENTPSLIST(PARALTLIST), FALSE);
4286                  WHILE CHECKDELIM( [GRPAREN]) = INVALIDTOKEN DO ERROR(14,0);
4287                  GATEPS := NIL
4288                  END
4289
4290      END; { of case }
4291
4292      IF GATEPS <> NIL THEN SETCOMBARS(PARALTLIST, GATEPS);
4293      IF DELIMCHECK <> INVALIDTOKEN THEN NEXTTOKEN; { TRANSLATENOMEN has taken NEXTTOKEN }
4294      WHILE CHECKDELIM( [GLSQUARE]+TERMINATORS)=INVALIDTOKEN DO ERROR(24,0);
4295      IF TOKEN.DELIMVAL = GLSQUARE THEN MODIFYCHILDPOSITIONS(PARALTLIST);
4296      WHILE CHECKDELIM(TERMINATORS)=INVALIDTOKEN DO ERROR(24,0)
4297      END;
4298
4299      { of Procedure ELEMENT
-----}
4300
4301
4302
4303      {-----}
4304      PROCEDURE ALTNTVE(VAR PARALTLIST : PPALTBARS;
4305                      OPTIONALSUB   : BOOLEAN);
4306
4307      { Called once for each alternative in a definition expresssion
4308        and cycles round any number of levels of further substitution. Beyond the
4309        very first level, the bond type attachment is NOTSPECIFIED until the child
4310        structure itself is reached (in ELEMENT).
4311        At the end of each cycle, the former PARALTLIST
4312        is destroyed as a new one is created. The last one remaining is passed back
4313        (via the VAR parameter) to ALTNVLIST, where it is destroyed.

```

```

4314         Called by ALTNVLIST}
4315
4316     VAR DELIMCHECK      : DELIMTYPE;
4317         COMB            : PCOMBINLIST;
4318         NEWPARALTLIST,
4319         PTR              : PPALTBARS;
4320
4321
4322
4323     FUNCTION ALREADYINLIST(PTRPS      : PTRPSTYPE;
4324                             PARALTLIST : PPALTBARS) : BOOLEAN;
4325
4326     { Returns TRUE if PTRPS is the PARSTRUCT field of any item in PARALTLIST.
4327       Called by ADDCOMBINPSS}
4328
4329     VAR PSFOUND : BOOLEAN;
4330
4331     BEGIN
4332     PSFOUND := FALSE;
4333     WHILE (PARALTLIST <> NIL) AND NOT PSFOUND DO WITH PARALTLIST^ DO
4334     BEGIN
4335         PSFOUND := PARSTRUCT = PTRPS;
4336         PARALTLIST := NEXT
4337     END;
4338     ALREADYINLIST := PSFOUND
4339     END;
4340
4341
4342
4343     PROCEDURE ADDPARALT(VAR NEWPARALTLIST : PPALTBARS;
4344                         NEWPARENT        : PTRPSTYPE);
4345
4346     { Adds one new item to NEWPARALTLIST, setting the ALTBAR field to NIL (because,
4347       as SB and OSB have a higher precedence than /, there cannot be any
4348       alternatives for this level)
4349       Called by ADDCOMBINPSS}
4350
4351     VAR NEWPA : PPALTBARS;
4352
4353     BEGIN

```

```

4354 NEW(NEWPA);
4355 WITH NEWPA^ DO
4356 BEGIN
4357     PARSTRUCT := NEWPARENT;
4358     ALTBAR := NIL;
4359     CONNBONDS.CONNECTIONS := NOTSET;
4360     COPYCHILDPS := FALSE;
4361     NEW(PRNTPOSNS);
4362     ECTRSIZE := ECTRSIZE + 9;
4363     PRNTPOSNS^.COMBINED := FALSE;
4364     GETAVAILABLEPOSITIONS(PARSTRUCT, PRNTPOSNS^.MEMBERS, 1);
4365     NEXT := NEWPARALTLIST
4366 END;
4367 NEWPARALTLIST := NEWPA
4368 END;
4369
4370
4371
4372 PROCEDURE ADDCOMBINPSS(COMBIN          : PCOMBINLIST;
4373                       VAR NEWPARALTLIST : PPALTBARS);
4374
4375 { Adds to NEWPARALTLIST the partial structures pointed to by the combination
4376   in COMBIN. If COMBIN is non-BOTTOMBAR then the procedure calls itself recursively
4377   for each of the ALTERNATIVES. A check is made (FUNCTION ALREADYINLIST) to
4378   prevent the same PS being included in the list more than once.
4379   If COMBIN is BOTTOMBAR, NEWPARALT is used to insert the new item in the
4380   list.
4381   Called by Body of ALTNTVE}
4382
4383 VAR ALTERN      : PALTERNLIST;
4384
4385 BEGIN
4386 IF COMBIN^.BOTTOMBAR
4387 THEN IF ALREADYINLIST(COMBIN^.CHILDPS, NEWPARALTLIST)
4388 THEN {don't duplicate}
4389 ELSE ADDPARALT(NEWPARALTLIST, COMBIN^.CHILDPS)
4390 ELSE BEGIN
4391     ALTERN := COMBIN^.ALTERNATIVES;
4392     REPEAT
4393     ADDCOMBINPSS(ALTERN^.COMBINATION, NEWPARALTLIST);

```

```

4394         ALTERN := ALTERN^.NEXT
4395         UNTIL ALTERN = NIL
4396         END
4397     END;
4398
4399
4400
4401     BEGIN {Body of Procedure ALTNTVE}
4402     REPEAT
4403     REPEAT ELEMENT(PARALTLIST, OPTIONALSUB)
4404     UNTIL CHECKDELIM([GAMPERSAND])=INVALIDTOKEN;
4405     DELIMCHECK := CHECKDELIM([GSB,GOSB]);
4406     IF DELIMCHECK <> INVALIDTOKEN THEN
4407     BEGIN
4408     NEWPARALTLIST := NIL;
4409     WHILE PARALTLIST <> NIL DO
4410     BEGIN
4411     IF PARALTLIST^.ALTBAR = NIL
4412     THEN COMB := PARALTLIST^.PARSTRUCT^.CHILDGATE
4413     ELSE COMB := PARALTLIST^.ALTBAR^.COMBINATION;
4414     WHILE COMB <> NIL DO
4415     BEGIN
4416     ADDCOMBINPSS(COMB, NEWPARALTLIST);
4417     COMB := COMB^.NEXT
4418     END;
4419     PTR := PARALTLIST^.NEXT;
4420     DISPOSE(PARALTLIST);
4421     PARALTLIST := PTR
4422     END;
4423     PARALTLIST := NEWPARALTLIST;
4424     OPTIONALSUB := DELIMCHECK = GOSB
4425     END
4426     UNTIL DELIMCHECK=INVALIDTOKEN
4427     END;
4428         { of Procedure ALTNTVE
4429     -----}
4430
4431
4432
4433     {-----}

```

```

4434 FUNCTION PPOSNS(DUMMYCOMBIN : PCOMBINLIST;
4435                 NEWPARENT   : PTRPSTYPE;
4436                 DUMMYSUBST  : SUBSTITUENT) : PTGROUPMEMS;
4437
4438 { Returns an appropriate parent positions set for NEWPARENT, after considering
4439   the positions available in it for the appropriate MAGNITUDE, and those
4440   positions specified for substitution in DUMMYCOMBIN. Outputs appropriate
4441   error messages if incompatibilities are detected.
4442   Called by ADDFURTHERSUBTN}
4443
4444 VAR AVAILPOSNS : ARRAY[1..3] OF INTEGSET;
4445     POSNS      : PTGROUPMEMS;
4446     LMAG       : TBOND MAG;
4447
4448
4449
4450 FUNCTION FINDBOTTOM(ALTERN      : PALTERNLIST;
4451                    FUNCTION MAG : TBOND MAG) : TBOND MAG;
4452
4453 { Traces down all the alternatives in the list headed by ALTERN using the formal
4454   function MAG. Returns the largest bond magnitude.
4455   Called by LMAGNOCHECKS
4456   LMAGCHECKS}
4457
4458 VAR LBOND,
4459     MBOND : TBOND MAG;
4460     COMB  : PCOMBINLIST;
4461
4462 BEGIN
4463 LBOND := 0;
4464 WHILE ALTERN <> NIL DO
4465   BEGIN
4466     COMB := ALTERN^.COMBINATION;
4467     WHILE COMB <> NIL DO
4468       BEGIN
4469         MBOND := MAG(COMB);
4470         IF MBOND > LBOND THEN LBOND := MBOND;
4471         COMB := COMB^.NEXT
4472       END;
4473     ALTERN := ALTERN^.NEXT

```

```

4474         END;
4475     FINDBOTTOM := LBOND
4476     END;
4477
4478
4479
4480     FUNCTION LMAGNOCHECKS(COMBIN : PCOMBINLIST) : TBOND MAG;
4481
4482     { Traces down COMBIN by calling itself recursively via FINDBOTOM.
4483       At the bottom bar, returns the largest bond MAGNITUDE.
4484       Called by LMAGCHECKS
4485         FINDBOTTOM (as formal parameter) }
4486
4487     VAR LMAG : TBOND MAG;
4488
4489     BEGIN
4490     IF COMBIN^.BOTTOMBAR
4491     THEN WITH COMBIN^, CONNBONDS DO
4492     CASE CONNECTIONS OF
4493     NOTSET : LMAG := 1;
4494     0      : LMAG := 0;
4495     1      : LMAG := MAGNITUDE(BOND);
4496     2      : BEGIN
4497             LMAG := MAGNITUDE(BONDA);
4498             IF MAGNITUDE(BONDB) > LMAG
4499             THEN LMAG := MAGNITUDE(BONDB)
4500             END
4501     END
4502     ELSE LMAG := FINDBOTTOM(COMBIN^.ALTERNATIVES, LMAGNOCHECKS);
4503     LMAGNOCHECKS := LMAG
4504     END;
4505
4506
4507
4508     FUNCTION LMAGCHECKS(COMBIN : PCOMBINLIST) : TBOND MAG;
4509
4510     { Traces down COMBIN, calling itself recursively via FINDBOTTOM,
4511       until it encounters a bar containing a PARENTPOSITIONS field,
4512       or reaches the bottom of the gate. In either case, calls LMAGNOCHECKS
4513       to obtain the largest bond magnitude from the

```

```

4514         bottom of the gate, and ensures that the PARENTPOSITIONS field (if given) is
4515         a subset of the positions available for this magnitude. The value of the
4516         magnitude is returned.
4517         Called by PPOSNS
4518             FINDBOTTOM (as formal parameter)}
4519
4520     VAR LMAG : TBOND MAG;
4521
4522     BEGIN
4523     IF COMBIN^.PARENTPOSITIONS <> NIL
4524     THEN BEGIN
4525         LMAG := LMAGNOCHECKS (COMBIN);
4526         WITH COMBIN^, PARENTPOSITIONS^ DO
4527             IF COMBINED
4528             THEN CHECKALLWITHIN (COMB MEMS, AVAIL POSNS [LMAG], 3)
4529             ELSE IF MEMBERS <= AVAIL POSNS [LMAG]
4530             THEN { specified positions are all available }
4531             ELSE FAILURE (49, DUMMYSUBST, ' ');
4532     END
4533     ELSE IF COMBIN^.BOTTOMBAR
4534     THEN LMAG := LMAGNOCHECKS (COMBIN)
4535     ELSE LMAG := FINDBOTTOM (COMBIN^.ALTERNATIVES, LMAGCHECKS);
4536     LMAGCHECKS := LMAG
4537     END;
4538
4539
4540
4541     BEGIN {Body of PPOSNS}
4542     FOR LMAG := 1 TO 3 DO GETAVAILABLEPOSITIONS (NEWPARENT, AVAIL POSNS [LMAG], LMAG);
4543     IF AVAIL POSNS [1] = []
4544     THEN FAILURE (50, DUMMYSUBST, ' ');
4545     WITH DUMMYCOMBIN^ DO
4546     IF PARENTPOSITIONS^.COMBINED
4547     THEN BEGIN
4548         LMAG := LMAGCHECKS (DUMMYCOMBIN);
4549         PPOSNS := PARENTPOSITIONS
4550     END
4551     ELSE IF PARENTPOSITIONS^.MEMBERS = [1..MAXCT]
4552     THEN BEGIN
4553         IF NOT DUMMYCOMBIN^.BOTTOMBAR

```

```

4554             THEN LMAG := FINDBOTTOM(DUMMYCOMBIN^.ALTERNATIVES, LMAGCHECKS);
4555             NEW(POSNS);
4556             ECTRSIZE := ECTRSIZE + 9;
4557             POSNS^.COMBINED := FALSE;
4558             POSNS^.MEMBERS := AVAILPOSNS[LMAG];
4559             PPOSNS := POSNS
4560             END
4561         ELSE BEGIN
4562             LMAG := LMAGCHECKS(DUMMYCOMBIN);
4563             PPOSNS := PARENTPOSITIONS
4564             END
4565     END;
4566             { of PPOSNS
4567     -----}
4568
4569
4570
4571     PROCEDURE ADDFURTHERSUBTN(COMBINBAR : PCOMBINLIST;
4572                             DUMMYPS   : PTRPSTYPE);
4573
4574     { Copies the further substitution on DUMMYPS onto the PSs at the bottom of
4575     COMBINBAR.
4576     Called by Body of ALTNVLIST}
4577
4578     VAR ALTERNBAR : PALTERNLIST;
4579         FSUBCOMB  : PCOMBINLIST;
4580
4581     BEGIN
4582     WHILE COMBINBAR <> NIL DO WITH COMBINBAR^ DO
4583     BEGIN
4584         IF BOTTOMBAR
4585         THEN BEGIN
4586             FSUBCOMB := DUMMYPS^.CHILDGATE;
4587             REPEAT
4588                 COPYCOMBAR({newcombar}    CHILDP^CHILDGATE,
4589                             {oldcombar}    FSUBCOMB,
4590                             {lastcomblayer} NIL,
4591                             {lastpposns}  PPOSNS(FSUBCOMB, CHILDP, DUMMYPS^.SUBSTNAME),
4592                             {prntps}      CHILDP,
4593                             {firstbar}    TRUE,

```

```

4594             {omitpg}           DEFNTABLEENTRY(CHILDPS),
4595             {copyyps}           FALSE);
4596             FSUBCOMB := FSUBCOMB^.NEXT
4597             UNTIL FSUBCOMB = NIL
4598             END
4599             ELSE BEGIN
4600                 ALTERNBAR := ALTERNATIVES;
4601                 REPEAT
4602                     ADDFURTHERSUBTN(ALTERNBAR^.COMBINATION, DUMMYPS);
4603                     ALTERNBAR := ALTERNBAR^.NEXT
4604                     UNTIL ALTERNBAR = NIL
4605                 END;
4606             COMBINBAR := NEXT
4607             END
4608         END;
4609
4610
4611
4612         BEGIN { Body of Procedure ALTNVLIST }
4613         REPEAT
4614             PARALTLIST := NIL;
4615             READPTR := PARENTPSLIST;
4616             WHILE READPTR <> NIL DO
4617                 BEGIN
4618                     NEW(NEWALTERNATIVE);
4619                     ECTRSIZE := ECTRSIZE + 4;
4620                     NEWALTERNATIVE^.COMBINATION := NIL;
4621                     NEWALTERNATIVE^.NEXT := READPTR^.COMBINS^.ALTERNATIVES;
4622                     READPTR^.COMBINS^.ALTERNATIVES := NEWALTERNATIVE;
4623
4624                     NEW(WRITEPTR);
4625                     WRITEPTR^.PARSTRUCT := READPTR^.PARSTRUCT;
4626                     WRITEPTR^.CONNBONDS := READPTR^.CONNBONDS;
4627                     WRITEPTR^.PRNTPOSNS := READPTR^.PRNTPOSNS;
4628                     WRITEPTR^.ALTBAR := NEWALTERNATIVE;
4629                     IF READPTR^.COPYCHILDPS
4630                         THEN WRITEPTR^.COPYCHILDPS := TRUE
4631                         ELSE IF READPTR^.FURTHERSUB = NIL
4632                             THEN WRITEPTR^.COPYCHILDPS := FALSE
4633                             ELSE WRITEPTR^.COPYCHILDPS := READPTR^.FURTHERSUB^.CHILDGATE <> NIL;

```

```

4634         WRITEPTR^.NEXT := PARALTLIST;
4635         PARALTLIST := WRITEPTR;
4636
4637         READPTR := READPTR^.NEXT
4638     END;
4639     ALTNVLIST(PARALTLIST, OPTIONALSUB);
4640     WHILE PARALTLIST <> NIL DO
4641     BEGIN
4642         WRITEPTR := PARALTLIST^.NEXT;
4643         DISPOSE(PARALTLIST);
4644         PARALTLIST := WRITEPTR
4645     END;
4646
4647     READPTR := PARENTPSLIST;
4648     WHILE READPTR <> NIL DO WITH READPTR^ DO
4649     BEGIN
4650         IF FURTHERSUB <> NIL
4651             THEN IF FURTHERSUB^.CHILDGATE <> NIL
4652                 THEN ADDFURTHERSUBTN(COMBINS^.ALTERNATIVES^.COMBINATION, FURTHERSUB);
4653         READPTR := NEXT
4654     END
4655
4656     UNTIL CHECKDELIM([GSLASH])=INVALIDTOKEN;
4657     WHILE PARENTPSLIST <> NIL DO
4658     BEGIN
4659         READPTR := PARENTPSLIST^.NEXT;
4660         DISPOSE(PARENTPSLIST);
4661         PARENTPSLIST := READPTR
4662     END
4663 END;
4664
4665                                     { OF PROCEDURE ALTNVLIST
4666     *****}
4667
4668
4669
4670
4671
4672
4673     {*****}

```

```

4674
4675                                PROCEDURE ASSIGNMENTSTMNT
4676
4677    *****}
4678    PROCEDURE ASSIGNMENTSTMNT;
4679
4680        { Analyses an assignment statement.
4681          Called by STATEMENT}
4682
4683    VAR SELECTEDFREQ : INTRECORD;
4684
4685
4686
4687    FUNCTION ASSGNTOP : TSELECTMODE;
4688
4689        { Returns the value of an assignment operator.
4690          Called by SUBSTASSIGNMENT
4691            MULTASSIGNMENT}
4692
4693    BEGIN
4694    WHILE CHECKDELIM([GDEQ,GSEQ,GHASHEQ,GDOLEQ,GEQUALS]) = INVALIDTOKEN DO ERROR(32,0);
4695    CASE TOKEN.DELIMVAL OF
4696        GEQUALS : ASSGNTOP := INDEPENDENT;
4697        GSEQ     : ASSGNTOP := ALLSAME;
4698        GDEQ     : ASSGNTOP := ALLDIFF;
4699        GDOLEQ   : ASSGNTOP := NOTALLSAME;
4700        GHASHEQ : ASSGNTOP := NOTALLDIFF
4701    END
4702    END;
4703
4704
4705
4706    {-----}
4707    PROCEDURE SUBSTGROUP(VAR GROUPMEMS : TGROUPMEMS;
4708                        LIMITSET      : INTEGSET;
4709                        ERRORCODE     : INTEGER);
4710
4711        { Analyses a substituent group, all members of which must fall in LIMITSET.
4712          Called by SUBSTASSIGNMENT}
4713

```

```
4714 LABEL 10;
4715
4716 VAR TERMINATORS : DELIMSET;
4717     CONNECTIVITY : TCONNS;
4718
4719
4720
4721 PROCEDURE REVISELIMITS(VAR LIMITSET : INTEGSET;
4722     CONNECTIVITY : TCONNS);
4723
4724 { Removes those elements in LIMITSET for which the connectivity shown in
4725   RDECLARATIONTABLE is not compatible with CONNECTIVITY}
4726
4727 VAR SUBST : SUBSTITUENT;
4728
4729 BEGIN
4730   FOR SUBST := 1 TO MAXVARS DO
4731     IF SUBST IN LIMITSET
4732       THEN WITH RDECLARATIONTABLE[SUBST]^ .CONN BONDS DO
4733         IF (CONNECTIONS=CONNECTIVITY) OR (CONNECTIONS=NOTSET)
4734           THEN {would be compatible}
4735           ELSE LIMITSET := LIMITSET - [SUBST]
4736   END;
4737
4738
4739
4740 PROCEDURE CHECKCOMPATABILITY(GROUP      : INTEGSET;
4741     CONNECTIVITY : TCONNS);
4742
4743 { Checks that the substituents in GROUP have compatible CONNECTIVITY }
4744
4745 VAR SUBST : SUBSTITUENT;
4746
4747 BEGIN
4748   IF CONNECTIVITY = NOTSET
4749     THEN FOR SUBST := 1 TO MAXVARS DO
4750       IF SUBST IN GROUP
4751         THEN WITH RDECLARATIONTABLE[SUBST]^ , CONN BONDS DO
4752           IF CONNECTIONS <> NOTSET
4753             THEN CONNECTIVITY := CONNECTIONS;
```

```
4754 IF CONNECTIVITY = NOTSET
4755 THEN {still no information on connectivity. No compatability checking possible}
4756 ELSE FOR SUBST := 1 TO MAXVARS DO
4757     IF SUBST IN GROUP
4758         THEN WITH RDECLARATIONTABLE[SUBST]^ DO
4759             IF CONNBONDS.CONNECTIONS = NOTSET
4760                 THEN BEGIN
4761                     SETCONN BONDS(CONNBONDS, CONNECTIVITY);
4762                     UPDATEPPSCONNS(RDECLARATIONTABLE[SUBST])
4763                 END
4764             ELSE IF CONNBONDS.CONNECTIONS = CONNECTIVITY
4765                 THEN {compatible}
4766                 ELSE FAILURE(51, 0, ' ')
4767 END;
4768
4769
4770 FUNCTION LISTSMATCH(SUBSTA, SUBSTB : SUBSTITUENT) : BOOLEAN;
4771
4772 { Checks that the declarations for SUBSTA and SUBSTB all refer to the same PSs }
4773
4774 VAR ADECNS,
4775     BDECNS : PPSLIST;
4776     FOUND : BOOLEAN;
4777
4778 BEGIN
4779 ADECNS := RDECLARATIONTABLE[SUBSTA];
4780 REPEAT
4781     BDECNS := RDECLARATIONTABLE[SUBSTB];
4782     REPEAT
4783         FOUND := (ADECNS^.PARSTRUCT = BDECNS^.PARSTRUCT) AND (ADECNS <> BDECNS);
4784         BDECNS := BDECNS^.NEXT
4785     UNTIL FOUND OR (BDECNS = NIL);
4786     IF FOUND
4787     THEN ADECNS := ADECNS^.NEXT
4788     ELSE BEGIN
4789         ERROR(37,0);
4790         ADECNS := NIL
4791     END
4792 UNTIL ADECNS = NIL;
4793
```

```

4794     LISTSMATCH := FOUND
4795     END;
4796
4797
4798
4799     PROCEDURE SUBSTCOMBINATION(VAR COMBMEMS : PDOUBLIST);
4800
4801     { Analyses a substituent combination }
4802
4803     VAR SUBCOMB : PDOUBLIST;
4804
4805     BEGIN
4806     NEW(SUBCOMB);
4807     NEXTTOKEN;
4808     WHILE CHECKDELIM([GR]) = INVALIDTOKEN DO ERROR(26,0);
4809     NEXTTOKEN;
4810     CHECKVALIDINT(LIMITSET,35);
4811     SUBCOMB^.FIRST := TOKEN.INTEGVAL;
4812     NEXTTOKEN;
4813     WHILE CHECKDELIM([GPLUS]) = INVALIDTOKEN DO ERROR(36, 0);
4814     NEXTTOKEN;
4815     WHILE CHECKDELIM([GR]) = INVALIDTOKEN DO ERROR(26, 0);
4816     NEXTTOKEN;
4817     WITH SUBCOMB^ DO
4818     REPEAT
4819     CHECKVALIDINT(LIMITSET, 35);
4820     SECOND := TOKEN.INTEGVAL
4821     UNTIL LISTSMATCH(FIRST, SECOND) AND LISTSMATCH(SECOND, FIRST);
4822     NEXTTOKEN;
4823     WHILE CHECKDELIM([GCOMMA] + TERMINATORS) = INVALIDTOKEN DO ERROR(24, 0);
4824     SUBCOMB^.NEXT := COMBMEMS;
4825     COMBMEMS := SUBCOMB
4826     END;
4827
4828
4829
4830     BEGIN { Body of Procedure SUBSTGROUP }
4831     TERMINATORS := [GNOTEQ,GEQUALS,GDEQ, GSEQ, GDOLEQ, GHASHEQ, GPERIOD];
4832     LOOKAHEAD;
4833     CHECKVALIDINT(LIMITSET, ERRORCODE);

```

```

4834 CONNECTIVITY := RDECLARATIONTABLE[TOKEN.INTEGVAL]^ .CONN BONDS.CONNECTIONS;
4835 LOOKAHEAD;
4836 10:
4837 WHILE CHECKDELIM([GPLUS, GCOMMA, GHYPHEN]+TERMINATORS)=INVALIDTOKEN DO ERROR(24,0);
4838 GROUPMEMS.COMBINED := TOKEN.DELIMVAL = GPLUS;
4839 IF GROUPMEMS.COMBINED
4840     THEN CASE CONNECTIVITY OF
4841         0, 2 : BEGIN
4842             ERROR(30,0);
4843             GOTO 10
4844         END;
4845     NOTSET,
4846     1 : BEGIN
4847         REVISELIMITS(LIMITSET, 1);
4848         GROUPMEMS.COMBMEMS := NIL;
4849         REPEAT SUBSTCOMBINATION(GROUPMEMS.COMBMEMS)
4850         UNTIL CHECKDELIM([GCOMMA])=INVALIDTOKEN
4851         END
4852     END
4853     ELSE BEGIN
4854         IF CONNECTIVITY <> NOTSET THEN REVISELIMITS(LIMITSET, CONNECTIVITY);
4855         NEXTTOKEN;
4856         GROUPRANGE(GROUPMEMS.MEMBERS, LIMITSET, ERRORCODE);
4857         CHECKCOMPATABILITY(GROUPMEMS.MEMBERS, CONNECTIVITY)
4858         END;
4859 WHILE CHECKDELIM(TERMINATORS)=INVALIDTOKEN DO ERROR(24,0)
4860 END;
4861
4862 { of Procedure SUBSTGROUP
4863 -----}
4864
4865
4866
4867 {-----}
4868 FUNCTION POINTERLIST(GROUPMEMS : TGROUPMEMS) : PPSLIST;
4869
4870 VAR READPTR,
4871     WRITEPTR,
4872     LISTPTR : PPSLIST;
4873     SUBST : SUBSTITUENT;

```

```
4874
4875 { Sets up a linked list of declarations for use in gate-setting.
4876   One element in the list represents one declaration of one substituent,
4877   so the same parent may appear several times in the list.
4878   The list is built up from the entries in RDECLARATIONTABLE.
4879   Called by SUBSTASSIGNMENT}
4880
4881
4882
4883 PROCEDURE GETBINFO(VAR BONDB : BONDORDER;
4884                   VAR POSNSB : INTEGSET;
4885                   BPTR      : PPSLIST;
4886                   PARENTPS  : PTRPSTYPE);
4887
4888 { Obtains bond order and positions from the second substituent of a combination.
4889   They are compiled from the information given in all the items in the BPTR
4890   list which reference PARENTPS.
4891   Called by ADDCOMBSUBS}
4892
4893 VAR FAILSTRING : STRING4;
4894
4895 BEGIN
4896   BONDB := NOTSPECIFIED;
4897   POSNSB := [];
4898   WHILE BPTR <> NIL DO WITH BPTR^ DO
4899     BEGIN
4900       IF PARSTRUCT = PARENTPS THEN
4901         BEGIN
4902           IF PRNTPOSNS^.COMBINED
4903             THEN PROGERROR(22) {COMBINED position set in combined substituent}
4904             ELSE POSNSB := POSNSB + PRNTPOSNS^.MEMBERS;
4905           CASE CONNBONDS.CONNECTIONS OF
4906             NOTSET : ;
4907             0, 2   : PROGERROR(26); {Only connectivity of 1 permitted in combined substituents}
4908             1     : IF CONNBONDS.BOND <> NOTSPECIFIED
4909                 THEN BEGIN
4910                   BONDB := BONDMATCHARRAY[CONNBONDS.BOND, BONDB];
4911                   IF BONDB = NOTSPECIFIED
4912                     THEN BEGIN
4913                     FAILSTRING[1] := BONDSTRING[CONNBONDS.BOND, 1];
```

```

4914                                FAILSTRING[2] := BONDSTRING[CONNBONDS.BOND, 2];
4915                                FAILSTRING[3] := BONDSTRING[BONDB, 1];
4916                                FAILSTRING[4] := BONDSTRING[BONDB, 2];
4917                                FAILURE(42, 0, FAILSTRING)
4918                                END
4919                                END
4920                                END
4921                                END;
4922                                BPTR := NEXT
4923                                END;
4924                                IF POSNSB = [] THEN PROGERROR(25)
4925                                END;
4926
4927
4928                                PROCEDURE ADDCOMBSUBS(APTR,
4929                                                                BPTR          : PPSLIST;
4930                                                                VAR LISTPTR : PPSLIST);
4931
4932                                { Adds a substituent combination to the PPSLIST.
4933                                  Called by Body of POINTERLIST}
4934
4935                                VAR WRITEPTR : PPSLIST;
4936                                    POSNSA,
4937                                    POSNSB,
4938                                    POSNSC   : INTEGSET;
4939                                    MAGSUM   : INTEGER;
4940
4941                                BEGIN
4942                                NEW(WRITEPTR);
4943                                WITH WRITEPTR^ DO
4944                                BEGIN
4945                                    PARSTRUCT := APTR^.PARSTRUCT;
4946                                    FURTHERSUB := NIL;
4947                                    WITH CONNBONDS DO
4948                                    BEGIN
4949                                        CONNECTIONS := 2;
4950                                        CASE APTR^.CONNBONDS.CONNECTIONS OF
4951                                        NOTSET : BONDA := NOTSPECIFIED;
4952                                        0, 2   : PROGERROR(23); {connectivity must be 1 for combined substituents}
4953                                        1     : BONDA := APTR^.CONNBONDS.BOND

```

```

4954         END;
4955         GETBINFO(BONDB, POSNSB, BPTR, PARSTRUCT)
4956     END;
4957     WITH APTR^.PRNTPOSNS^ DO
4958         IF COMBINED
4959             THEN PROGERROR(24) {COMBINED position set with combined substituents}
4960             ELSE POSNSA := MEMBERS;
4961     IF POSNSA * POSNSB = []
4962     THEN POSNSC := []
4963     ELSE BEGIN
4964         WITH CONNBONDS DO MAGSUM := MAGNITUDE(BONDA) + MAGNITUDE(BONDB);
4965         IF MAGSUM <= 3
4966             THEN GETAVAILABLEPOSITIONS(PARSTRUCT, POSNSC, MAGSUM)
4967             ELSE POSNSC := [];
4968         POSNSC := POSNSA * POSNSB * POSNSC
4969     END;
4970     NEW(PRNTPOSNS);
4971     ECTRSIZE := ECTRSIZE + 9;
4972     WITH PRNTPOSNS^ DO
4973         BEGIN
4974             COMBINED := TRUE;
4975             COMMEMS := NIL;
4976             LISTPOSNS(COMMEMS, POSNSA, POSNSB, POSNSC)
4977         END;
4978     NEW(COMBINS, TRUE);
4979     ECTRSIZE := ECTRSIZE + 11;
4980     WITH COMBINS^ DO
4981         BEGIN
4982             PARENTPOSITIONS := PRNTPOSNS;
4983             FREQUENCY.TOPRANGE := NOTSET;
4984             FREQUENCY.SUBRANGES := ESSENTFREQ;
4985             BOTTOMBAR := FALSE;
4986             ALTERNATIVES := NIL;
4987             NEXT := PARSTRUCT^.CHILDGATE
4988         END;
4989     PARSTRUCT^.CHILDGATE := WRITEPTR^.COMBINS;
4990     COPYCHILDPS := FALSE;
4991     NEXT := LISTPTR
4992 END;
4993 LISTPTR := WRITEPTR

```

```

4994     END;
4995
4996
4997
4998     PROCEDURE ADDDEFNTABLE(VAR LISTPTR : PPSLIST;
4999                           SUBST      : SUBSTITUENT);
5000
5001     { Adds the RDEFINITIONTABLE for SUBST to the PPSLIST in LISTPTR. If the SUBST
5002       has already been defined then the COMBINS is taken from RDEFINITIONTABLE,
5003       otherwise a new one is created and entered into RDEFINITIONTABLE.
5004       Called by Body of POINTERLIST}
5005
5006     VAR WRITEPTR : PPSLIST;
5007
5008     BEGIN
5009     NEW(WRITEPTR);
5010     WITH WRITEPTR^ DO
5011     BEGIN
5012     PARSTRUCT := NIL;
5013     FURTHERSUB := NIL;
5014     NEW(PRNTPOSNS);
5015     ECTRSIZE := ECTRSIZE + 9;
5016     PRNTPOSNS^.COMBINED := FALSE;
5017     PRNTPOSNS^.MEMBERS := [1..MAXCT];
5018     SETCONNBOARDS(CONNBONDS, LISTPTR^.CONNBONDS.CONNECTIONS);
5019     COPYCHILDPS := FALSE;
5020     NEXT := LISTPTR;
5021     IF RDEFINITIONTABLE[SUBST] = NIL
5022     THEN BEGIN
5023     NEW(COMBINS, TRUE);
5024     WITH COMBINS^ DO
5025     BEGIN
5026     BOTTOMBAR := FALSE;
5027     ALTERNATIVES := NIL;
5028     NEXT := NIL;
5029     PARENTPOSITIONS := PRNTPOSNS;
5030     FREQUENCY.TOPRANGE := NOTSET;
5031     FREQUENCY.SUBRANGES := NIL
5032     END;
5033     RDEFINITIONTABLE[SUBST] := COMBINS

```

```

5034         END
5035     ELSE COMBINS := RDEFINITIONTABLE[SUBST]
5036     END;
5037 LISTPTR := WRITEPTR
5038 END;
5039
5040
5041
5042 BEGIN {Body of POINTERLIST}
5043 LISTPTR := NIL;
5044 IF GROUPMEMS.COMBINED
5045     THEN WHILE GROUPMEMS.COMBMEMS <> NIL DO WITH GROUPMEMS DO
5046         BEGIN
5047             READPTR := RDECLARATIONTABLE[COMBMEMS^.FIRST];
5048             WHILE READPTR <> NIL DO
5049                 BEGIN
5050                     ADDCOMBSUBS(READPTR, RDECLARATIONTABLE[COMBMEMS^.SECOND], LISTPTR);
5051                     READPTR := READPTR^.NEXT
5052                 END;
5053                 COMBMEMS := COMBMEMS^.NEXT
5054             END
5055         ELSE FOR SUBST := 1 TO MAXVARS DO IF SUBST IN GROUPMEMS.MEMBERS
5056             THEN BEGIN
5057                 READPTR := RDECLARATIONTABLE[SUBST];
5058                 WHILE READPTR <> NIL DO
5059                     BEGIN
5060                         NEW(WRITEPTR);
5061                         WRITEPTR^. := READPTR^;
5062                         IF WRITEPTR^.COMBINS = NIL
5063                             THEN PROGERROR(21); {Declaration without combination bar}
5064                         WRITEPTR^.NEXT := LISTPTR;
5065                         LISTPTR := WRITEPTR;
5066                         READPTR := READPTR^.NEXT
5067                     END;
5068                     ADDDEFNTABLE(LISTPTR, SUBST)
5069                 END;
5070             POINTERLIST := LISTPTR
5071         END;
5072     {of FUNCTION POINTERLIST
-----}
5073

```

```

5074
5075
5076
5077 PROCEDURE SUBSTASSIGNMENT;
5078
5079 { Analyses a substituent assignment.
5080   Called by body of ASSIGNMENTSTMNT}
5081
5082 VAR GROUPMEMS      : TGROUPMEMS;
5083     PARENTPSLIST   : PPSLIST;
5084     DELPTR          : PDOUBLIST;
5085
5086 BEGIN
5087   SUBSTGROUP(GROUPMEMS, DECLSUBS, 1);
5088   IF ASSGNTOP <> INDEPENDENT THEN
5089     WRITELN('Non-independent assignment not yet implemented');
5090   PARENTPSLIST := POINTERLIST(GROUPMEMS);
5091   ALTNVLIST(PARENTPSLIST, FALSE);
5092   WITH GROUPMEMS DO IF COMBINED
5093     THEN WHILE COMBMEMS <> NIL DO
5094       BEGIN
5095         DEFNSUBS := DEFNSUBS + [COMBMEMS^.FIRST, COMBMEMS^.SECOND];
5096         DELPTR := COMBMEMS;
5097         COMBMEMS := COMBMEMS^.NEXT;
5098         DISPOSE(DELPTR)
5099       END
5100     ELSE DEFNSUBS := DEFNSUBS + MEMBERS
5101 END;
5102
5103
5104
5105 {-----}
5106 PROCEDURE MULTASSIGNMENT;
5107
5108 { Analyses a multiplier assignment.
5109   Called by Body of ASSIGNMENTSTMNT}
5110
5111 VAR MULT           : MULTIPLIER;
5112     DEFINEDMULTS   : INTEGSET;
5113     MULTVALUES,

```

```

5114         MULTVALCOPY : INTRECORD;
5115
5116
5117
5118     PROCEDURE COPYLIST(MULTVALUES      : INTRECORD;
5119                       VAR MULTVALCOPY : INTRECORD);
5120
5121     { Copies a list of values for a multiplier.
5122       Called by body of MULTASSIGNMENT}
5123
5124     VAR NEWITEM,
5125         LASTITEM : PDOUBLIST;
5126
5127     BEGIN
5128     MULTVALCOPY.TOPRANGE := MULTVALUES.TOPRANGE;
5129     MULTVALCOPY.SUBRANGES := NIL;
5130     WHILE MULTVALUES.SUBRANGES <> NIL DO
5131     BEGIN
5132     NEW(NEWITEM);
5133     ECTRSIZE := ECTRSIZE + 6;
5134     WITH NEWITEM^ DO
5135     BEGIN
5136     FIRST := MULTVALUES.SUBRANGES^.FIRST;
5137     SECOND := MULTVALUES.SUBRANGES^.SECOND;
5138     NEXT := NIL
5139     END;
5140     IF MULTVALCOPY.SUBRANGES = NIL
5141     THEN MULTVALCOPY.SUBRANGES := NEWITEM
5142     ELSE LASTITEM^.NEXT := NEWITEM;
5143     LASTITEM := NEWITEM;
5144     MULTVALUES.SUBRANGES := MULTVALUES.SUBRANGES^.NEXT
5145     END
5146     END;
5147
5148
5149
5150     PROCEDURE ADDITEM(VAR NEWITEM,
5151                     NEWLIST : PDOUBLIST);
5152
5153     { Inserts NEWITEM into NEWLIST.

```

```

5154         Called by COMBINEVALUES}
5155
5156     VAR NEWLISTITEM : PDOUBLIST;
5157
5158     BEGIN
5159     NEWLISTITEM := NEWITEM;
5160     NEWITEM := NEWITEM^.NEXT;
5161     NEWLISTITEM^.NEXT := NEWLIST;
5162     NEWLIST := NEWLISTITEM
5163     END;
5164
5165
5166
5167     PROCEDURE COMBINEVALUES(VAR TABLEVALUES : INTRECORD;
5168                             NEWVALUES      : INTRECORD);
5169
5170     { Combines the NEWVALUES just obtained with those already in TABLEVALUES.
5171       Called by Body of MULTASSIGNMENT }
5172
5173     VAR NEWLIST,
5174         NEWITEM : PDOUBLIST;
5175         FINISHED : BOOLEAN;
5176
5177     BEGIN
5178     IF (TABLEVALUES.TOPRANGE) = NOTSET
5179     THEN TABLEVALUES.TOPRANGE := NEWVALUES.TOPRANGE
5180     ELSE IF NEWVALUES.TOPRANGE = NOTSET
5181     THEN { leave TABLEVALUES.TOPRANGE as it is }
5182     ELSE IF NEWVALUES.TOPRANGE < TABLEVALUES.TOPRANGE
5183     THEN TABLEVALUES.TOPRANGE := NEWVALUES.TOPRANGE;
5184     NEWLIST := NIL;
5185     WHILE NOT ((TABLEVALUES.SUBRANGES = NIL) AND (NEWVALUES.SUBRANGES = NIL)) DO
5186     IF TABLEVALUES.SUBRANGES = NIL
5187     THEN ADDITEM(NEWVALUES.SUBRANGES, NEWLIST)
5188     ELSE IF NEWVALUES.SUBRANGES = NIL
5189     THEN ADDITEM(TABLEVALUES.SUBRANGES, NEWLIST)
5190     ELSE IF TABLEVALUES.SUBRANGES^.FIRST > NEWVALUES.SUBRANGES^.FIRST
5191     THEN ADDITEM(TABLEVALUES.SUBRANGES, NEWLIST)
5192     ELSE ADDITEM(NEWVALUES.SUBRANGES, NEWLIST);
5193     IF NEWLIST <> NIL

```

```

5194     THEN BEGIN
5195         TABLEVALUES.SUBRANGES := NEWLIST;
5196         NEWLIST := NEWLIST^.NEXT;
5197         TABLEVALUES.SUBRANGES^.NEXT := NIL
5198     END
5199     ELSE TABLEVALUES.SUBRANGES := NIL;
5200 WHILE NEWLIST <> NIL DO
5201     BEGIN
5202         NEWITEM := NEWLIST;
5203         NEWLIST := NEWLIST^.NEXT;
5204         IF NEWITEM^.FIRST > TABLEVALUES.SUBRANGES^.SECOND + 1
5205         THEN BEGIN
5206             NEWITEM^.NEXT := TABLEVALUES.SUBRANGES;
5207             TABLEVALUES.SUBRANGES := NEWITEM
5208         END
5209         ELSE BEGIN
5210             IF NEWITEM^.SECOND > TABLEVALUES.SUBRANGES^.SECOND
5211             THEN TABLEVALUES.SUBRANGES^.SECOND := NEWITEM^.SECOND;
5212             DISPOSE(NEWITEM);
5213             ECTRSIZE := ECTRSIZE - 6
5214         END
5215     END;
5216     FINISHED := (TABLEVALUES.SUBRANGES = NIL) OR (TABLEVALUES.TOPRANGE = NOTSET);
5217 WHILE NOT FINISHED DO WITH TABLEVALUES DO
5218     IF SUBRANGES^.SECOND >= TOPRANGE - 1
5219     THEN BEGIN
5220         IF SUBRANGES^.FIRST < TOPRANGE
5221         THEN TOPRANGE := SUBRANGES^.FIRST;
5222         NEWITEM := SUBRANGES;
5223         SUBRANGES := SUBRANGES^.NEXT;
5224         DISPOSE(NEWITEM);
5225         ECTRSIZE := ECTRSIZE - 6;
5226         FINISHED := (SUBRANGES = NIL)
5227     END
5228     ELSE FINISHED := TRUE
5229 END;
5230
5231
5232
5233 BEGIN {Body of MULTASSIGNMENT}

```

```

5234 GROUPRANGE(DEFINEDMULTS, DECLMULT, 2);
5235 WHILE CHECKDELIM([GDEQ,GSEQ,GHASHEQ,GDOLEQ,GNOTEQ,GEQUALS,GPERIOD]) = INVALIDTOKEN DO ERROR(24,0);
5236 IF ASSGNTOP <> INDEPENDENT
5237     THEN WRITELN('Non-independent assignment not yet implemented');
5238 NEXTTOKEN;
5239 SELECTOR(MULTVALUES,[0..MAXVARS],10);
5240 NEXTTOKEN;
5241 FOR MULT := 1 TO MAXVARS DO IF MULT IN DEFINEDMULTS THEN
5242     BEGIN
5243         COPYLIST(MULTVALUES, MULTVALCOPY);
5244         IF MULT IN DEFNMULT
5245             THEN COMBINEVALUES(MDEFINITIONTABLE[MULT], MULTVALCOPY)
5246             ELSE MDEFINITIONTABLE[MULT] := MULTVALCOPY
5247         END;
5248         REDUCEECTR(MULTVALUES.SUBRANGES);
5249         DESTROY(MULTVALUES.SUBRANGES);
5250         DEFNMULT := DEFNMULT + DEFINEDMULTS
5251     END;          {of MULTASSIGNMENT
5252     -----}
5253
5254
5255
5256 BEGIN { Body of ASSIGNMENTSTMNT }
5257 IF CHECKDELIM([GR,GM])=INVALIDTOKEN THEN SELECTOR(SELECTEDFREQ, [0..MAXVARS],9);
5258 WHILE CHECKDELIM([GR,GM]) = INVALIDTOKEN DO ERROR(20,0);
5259 IF TOKEN.DELIMVAL = GR
5260     THEN SUBSTASSIGNMENT
5261     ELSE MULTASSIGNMENT
5262 END;
5263     {of PROCEDURE ASSIGNMENTSTMNT
5264     *****}
5265
5266
5267
5268
5269 PROCEDURE CONDITION;
5270
5271 { Analyses a condition.
5272   Called by IFSTATEMENT
5273   RESTRICTSTMNT}

```

```
5274
5275 BEGIN
5276 CONDITIONSPRESENT := TRUE;
5277 REPEAT NEXTTOKEN
5278 UNTIL CHECKDELIM([GTHEN,GEND,GELSE,GSEMI,GPERIOD]) <> INVALIDTOKEN
5279 END;
5280
5281
5282
5283 PROCEDURE RESTRICTSTMNT;
5284
5285 { Analyses a RESTRICT statement.
5286   Called by STATEMENT}
5287
5288 BEGIN
5289 CONDITION;
5290 WHILE CHECKDELIM([GELSE,GEND,GSEMI,GPERIOD]) = INVALIDTOKEN DO ERROR(24,0)
5291 END;
5292
5293
5294
5295 PROCEDURE STATEMENT;
5296 FORWARD;
5297
5298
5299
5300 PROCEDURE IFSTATEMENT;
5301
5302 { Analyses an IF statement.
5303   Called by STATEMENT}
5304
5305 BEGIN
5306 CONDITION;
5307 WHILE CHECKDELIM ([GTHEN]) <> GTHEN DO ERROR(17,0);
5308 STATEMENT;
5309 WHILE CHECKDELIM ([GELSE,GEND,GSEMI,GPERIOD]) = INVALIDTOKEN DO ERROR(24,0);
5310 IF CHECKDELIM([GELSE]) = GELSE THEN
5311 BEGIN
5312 STATEMENT;
5313 WHILE CHECKDELIM ([GELSE,GEND,GSEMI,GPERIOD]) = INVALIDTOKEN DO ERROR(24,0)
```

```
5314         END
5315     END;
5316
5317
5318
5319     PROCEDURE CMPDSTMNT;
5320
5321     { Analyses a compound statement.
5322       Called by STATEMENT}
5323
5324     BEGIN
5325     REPEAT
5326         STATEMENT;
5327         WHILE CHECKDELIM([GSEMI,GEND]) = INVALIDTOKEN DO ERROR(15,0)
5328     UNTIL TOKEN.DELIMVAL = GEND;
5329     NEXTTOKEN
5330     END;
5331
5332
5333
5334     PROCEDURE STATEMENT;
5335
5336     { Analyses a STATEMENT, calling the appropriate procedure.
5337       Called by IFSTATEMENT
5338         RESTRICTSTMNT
5339         CMPDSTMNT
5340         Body of INTERPRET}
5341
5342     VAR DELIMCHECK : DELIMTYPE;
5343
5344     BEGIN
5345     NEXTTOKEN;
5346     REPEAT
5347         DELIMCHECK := CHECKDELIM([GBEGIN,GIF,GRESTRICT,GR,GM,GOPENANG,GEND,GSEMI,GPERIOD]);
5348         IF DELIMCHECK=INVALIDTOKEN THEN ERROR(19,0)
5349     UNTIL DELIMCHECK<>INVALIDTOKEN;
5350     CASE DELIMCHECK OF
5351         GR, GM,
5352         GOPENANG,
5353         INVALIDTOKEN : ASSIGNMENTSTMNT;
```

```

5354      GBEGIN      : CMPDSTMNT;
5355      GIF          : IFSTATEMENT;
5356      GRESTRIC    : RESTRICTSTMNT;
5357      GPERIOD,
5358      GEND, GSEMI  : (* empty statement *)
5359      END
5360      END;
5361
5362
5363
5364      PROCEDURE CHECKALLDONE;
5365
5366      { Checks that all declared substituents and multipliers have been defined
5367        Called by Body of INTERPRET}
5368
5369      VAR M          : 1..MAXVARS;
5370
5371      BEGIN
5372      IF (DECLSUBS - DEFNSUBS) <> [] THEN
5373          BEGIN
5374              WRITELN('The following substituents',' remain undefined:');
5375              FOR M := 1 TO MAXVARS DO
5376                  IF M IN (DECLSUBS-DEFNSUBS) THEN WRITE('  R', M:1);
5377              WRITELN;
5378              TOKEN.DELIMVAL := INVALIDTOKEN
5379          END;
5380      IF (DECLMULT - DEFNMULT) <> [] THEN
5381          BEGIN
5382              WRITELN('The following multipliers',' remain undefined:');
5383              FOR M := 1 TO MAXVARS DO
5384                  IF M IN (DECLMULT-DEFNMULT) THEN WRITE('  M', M:1);
5385              WRITELN;
5386              TOKEN.DELIMVAL := INVALIDTOKEN
5387          END;
5388      IF TOKEN.DELIMVAL = INVALIDTOKEN THEN
5389          WHILE CHECKDELIM([GSEMI])=INVALIDTOKEN DO ERROR(13,0)
5390      END;
5391
5392
5393

```

```
5394 PROCEDURE RECORDMULTS;
5395
5396 { Adds the values for multipliers to the appropriate FREQUENCY fields in the ECTR.
5397   Called by Body of INTERPRET}
5398
5399 VAR MULT : MULTIPLIER;
5400     PMPTR : PMDECLIST;
5401
5402 BEGIN
5403 FOR MULT := 1 TO MAXVARS DO IF MULT IN DEFNMULT THEN
5404     BEGIN
5405     PMPTR := MDECLARATIONTABLE[MULT];
5406     REPEAT
5407     PMPTR^.SUBSTDECN^.COMBINS^.FREQUENCY := MDEFINITIONTABLE[MULT];
5408     PMPTR := PMPTR^.NEXT
5409     UNTIL PMPTR = NIL
5410     END
5411 END;
5412
5413
5414
5415 {*****}
5416 PROCEDURE OUTINTREP;
5417
5418 { Outputs a representation of the ECTR to a diagnostics file }
5419
5420 VAR SUBST : SUBSTITUENT;
5421     MULT : MULTIPLIER;
5422     READPTR : PIRLIST;
5423
5424
5425
5426 FUNCTION PSNO(PTRPS : PTRPSTYPE) : INTEGER;
5427
5428 VAR PTR : PIRLIST;
5429     NUM : INTEGER;
5430     FOUND : BOOLEAN;
5431
5432 BEGIN
5433 NUM := 0;
```

```
5434 PTR := INTERNALREP.PSLIST;
5435 FOUND := FALSE;
5436 WHILE (PTR<>NIL) AND NOT FOUND DO
5437   BEGIN
5438     NUM := NUM + 1;
5439     IF PTR^.PARSTRUCT = PTRPS
5440       THEN FOUND := TRUE
5441       ELSE PTR := PTR^.NEXT
5442   END;
5443 IF FOUND THEN PSNO := NUM
5444   ELSE PSNO := 0
5445 END;
5446
5447
5448
5449 PROCEDURE WRITEFREQ(FREQUENCY : INTRECORD);
5450
5451 VAR PTR : PDOUBLIST;
5452
5453 BEGIN
5454 WRITE(DIAGFILE, '<');
5455 WITH FREQUENCY DO
5456   BEGIN
5457     IF TOPRANGE <> NOTSET THEN WRITE(DIAGFILE, TOPRANGE:1, '-',');
5458     PTR := SUBRANGES;
5459     WHILE PTR <> NIL DO WITH PTR^ DO
5460       BEGIN
5461         WRITE(DIAGFILE, FIRST:1, '-', SECOND:1, ',');
5462         PTR := NEXT
5463       END
5464     END;
5465 WRITE(DIAGFILE, '>')
5466 END;
5467
5468
5469
5470 PROCEDURE WRITEPOSNS(POSNSSET : TGROUPMEMS);
5471
5472 VAR POSN : ATOMNUMBER;
5473     PTR : PDOUBLIST;
```

```

5474
5475 BEGIN
5476 WRITE(DIAGFILE, '[');
5477 IF POSNSET.COMBINED
5478     THEN BEGIN
5479         PTR := POSNSET.COMBMEMS;
5480         WHILE PTR <> NIL DO WITH PTR^ DO
5481             BEGIN
5482                 WRITE(DIAGFILE, FIRST:1, '/', SECOND:1, ',');
5483                 PTR := NEXT;
5484             END
5485         END
5486     ELSE FOR POSN := 1 TO MAXCT DO IF POSN IN POSNSET.MEMBERS
5487         THEN WRITE(DIAGFILE, POSN:1, ',');
5488 WRITE(DIAGFILE, ']')
5489 END;
5490
5491
5492
5493 PROCEDURE WRITECONNS(CONNBONDS : TCONNBONDS);
5494
5495 { Writes out the bond orders in CONNBONDS.
5496   Called by WRITEPGS
5497     WRITECOMBIN
5498     WRITEDECN}
5499
5500 BEGIN
5501 WITH CONNBONDS DO
5502 CASE CONNECTIONS OF
5503     NOTSET : Writeln(DIAGFILE, '--');
5504     0      : ;
5505     1      : Writeln(DIAGFILE, BONDSTRING[BOND]);
5506     2      : BEGIN
5507         WRITE(DIAGFILE, BONDSTRING[BONDA]);
5508         Writeln(DIAGFILE, BONDSTRING[BONDB])
5509     END
5510 END
5511 END;
5512
5513

```

```

5514
5515 PROCEDURE WRITEPGS(PARENTS : PPARENTLIST);
5516
5517 { Writes a series of Parent Gates, headed by PARENTS, to DIAGFILE }
5518
5519
5520 BEGIN
5521 Writeln(DIAGFILE);
5522 Writeln(DIAGFILE, 'PARENT GATES: ');
5523 WHILE PARENTS <> NIL DO WITH PARENTS^ DO
5524 BEGIN
5525 WRITEPOSNS(CHILDPOSITIONS);
5526 WRITE(DIAGFILE, ' PS:', PSNO(PARENTPS) : 1, ' ');
5527 WRITEPOSNS(PARENTPOSITIONS);
5528 WRITECONNS(CONNBONDS);
5529 PARENTS := NEXT
5530 END
5531 END;
5532
5533
5534 PROCEDURE WRITECG(COMBINLIST : PCOMBINLIST;
5535 INDENT : INTEGER);
5536
5537 FORWARD;
5538
5539
5540 PROCEDURE WRITEALTERNs(ALTERNATIVES : PALTERNLIST;
5541 INDENT : INTEGER);
5542
5543
5544 VAR ALTNO,
5545 I : INTEGER;
5546
5547 BEGIN
5548 ALTNO := 0;
5549 WHILE ALTERNATIVES <> NIL DO
5550 BEGIN
5551 FOR I := 1 TO INDENT DO WRITE(DIAGFILE, ' ');
5552 ALTNO := ALTNO + 1;
5553 Writeln(DIAGFILE, ' ALT ', ALTNO : 2);

```

```

5554     WRITECG(ALTERNATIVES^.COMBINATION, INDENT + 4);
5555     ALTERNATIVES := ALTERNATIVES^.NEXT
5556     END
5557 END;
5558
5559
5560
5561     PROCEDURE WRITECOMBIN(COMBINPTR : PCOMBINLIST;
5562                          ITEMNO,
5563                          INDENT    : INTEGER);
5564
5565     { Outputs the information in the single combination gate pointed to by COMBINPTR. }
5566
5567     VAR I : INTEGER;
5568
5569     BEGIN
5570     WITH COMBINPTR^ DO
5571         BEGIN
5572             FOR I := 1 TO INDENT DO WRITE(DIAGFILE, ' ');
5573             WRITE(DIAGFILE, 'Item. No.', ITEMNO : 3 );
5574             IF COMBINPTR^.PARENTPOSITIONS = NIL
5575                 THEN WRITE(DIAGFILE, '[NIL]')
5576                 ELSE WRITEPOSNS(PARENTPOSITIONS^);
5577             WRITEFREQ(FREQUENCY);
5578             IF BOTTOBAR
5579                 THEN BEGIN
5580                     WRITE(DIAGFILE, 'PS:', PSNO(COMBINPTR^.CHILDPS) : 1);
5581                     WRITEPOSNS(CHILDPOSITIONS);
5582                     WRITECONNS(CONNBONDS)
5583                 END
5584                 ELSE BEGIN
5585                     WRITELN(DIAGFILE, 'ALTERNATIVES:');
5586                     WRITEALTERN(ALTERNATIVES, INDENT);
5587                     FOR I := 1 TO INDENT DO WRITE(DIAGFILE, ' ');
5588                     WRITELN(DIAGFILE, 'End of Item ', ITEMNO : 1, ' alternatives')
5589                 END
5590         END
5591     END;
5592
5593

```

```
5594
5595 PROCEDURE WRITECG; {FORWARD declaration above WRITEALTERNS}
5596
5597 VAR ITEMNO : INTEGER;
5598
5599 BEGIN
5600 ITEMNO := 0;
5601 WHILE COMBINLIST <> NIL DO
5602     BEGIN
5603         ITEMNO := ITEMNO + 1;
5604         WRITECOMBIN(COMBINLIST, ITEMNO, INDENT);
5605         COMBINLIST := COMBINLIST^.NEXT
5606     END
5607 END;
5608
5609
5610 PROCEDURE WRITECT(VAR CT : CCTYPE);
5611
5612 VAR ROWNO : ATOMNUMBER;
5613     CNGNR : 1..MAXCONGENERS;
5614
5615 BEGIN
5616 WRITELN(DIAGFILE, 'SPECIFIC': 10);
5617 FOR ROWNO := 1 TO MAXCT DO IF CT[ROWNO] <> NIL THEN WITH CT[ROWNO]^ DO
5618     BEGIN
5619         WRITE(DIAGFILE, ROWNO :2);
5620         IF ATOMICROW THEN WRITE(DIAGFILE, ATOM : 3)
5621             ELSE WRITE(DIAGFILE, 'R', NAME:1, ' ');
5622         WRITE(DIAGFILE, CHARGE : 3, HYDROGENS : 2);
5623         FOR CNGNR := 1 TO MAXCONGENERS DO WITH CONGENERS[CNGNR] DO
5624             IF RELATIONSHIP = FRATERNAL THEN
5625                 BEGIN
5626                     WRITE(DIAGFILE, ORD(BOND) : 3);
5627                     WRITE(DIAGFILE, ROWNUM : 3);
5628                 END;
5629             IF NOT ATOMICROW AND (VALUES <> NIL)
5630                 THEN BEGIN
5631                     WRITELN(DIAGFILE, 'VALUES:');
5632                     WRITECG(VALUES, 0);
5633
```

```

5634         END
5635     ELSE WRITELN(DIAGFILE)
5636     END
5637 END;
5638
5639
5640
5641 PROCEDURE WRITEPS(PTRPS      : PTRPSTYPE);
5642
5643
5644 BEGIN
5645 WRITE(DIAGFILE, PSNO(PTRPS):2 ,':');
5646 CASE PTRPS^.PSVARIETY OF
5647     DUMMY      : WRITELN(DIAGFILE,' DUMMY R', PTRPS^.SUBSTNAME: 1);
5648     UNKNOWN    : WRITELN(DIAGFILE,' UNKNOWN' );
5649     SPECIFIC   : WRITECT(PTRPS^.CT);
5650     GENERIC    : BEGIN
5651         WRITELN(DIAGFILE,' GENERIC');
5652         LISTPARAMS(DIAGFILE,PTRPS^.PARAMLIST)
5653     END;
5654     OTHER      : BEGIN
5655         WRITELN(DIAGFILE,'OTHER':7);
5656         WRITELN(DIAGFILE,PTRPS^.TERM)
5657     END
5658 END;
5659 IF PTRPS^.PARENTGATE <> NIL THEN WRITEPGS(PTRPS^.PARENTGATE);
5660 WRITELN(DIAGFILE);
5661 IF PTRPS^.CHILDGATE <> NIL THEN
5662     BEGIN
5663         WRITELN(DIAGFILE, 'CHILD GATES:');
5664         WRITECG(PTRPS^.CHILDGATE, 0)
5665     END;
5666 WRITELN(DIAGFILE, '-----');
5667 WRITELN(DIAGFILE)
5668 END;
5669
5670
5671
5672 PROCEDURE WRITEDECN(DECLPTR : PPSLIST);
5673

```

```

5674 { Writes out the information in DECLPTR. }
5675
5676 BEGIN
5677 WHILE DECLPTR <> NIL DO WITH DECLPTR^ DO
5678     BEGIN
5679         WRITELN(DIAGFILE);
5680         WRITELN(DIAGFILE, 'Declared in ', PSNO(PARSTRUCT) : 2);
5681         IF PRNTPOSNS = NIL THEN WRITE(DIAGFILE, '[NIL]')
5682             ELSE WRITEPOSNS(PRNTPOSNS^);
5683         WRITECONNS(CONNBONDS);
5684         IF FURTHERSUB <> NIL
5685             THEN WRITELN(DIAGFILE, 'Further substitution on PS ', PSNO(FURTHERSUB):2);
5686         IF COPYCHILDPS THEN WRITELN('COPYCHILDPS');
5687         DECLPTR := NEXT
5688     END
5689 END;
5690
5691
5692
5693 BEGIN { Body of OUTINTREP }
5694 WRITELN(DIAGFILE);
5695 WRITELN(DIAGFILE, '***** GENERIC STRUCTURE ', INTERNALREP.REFNUMBER : 3);
5696 WRITELN(DIAGFILE);
5697 WRITELN(DIAGFILE, 'Partial Structures: ');
5698 READPTR := INTERNALREP.PSLIST;
5699 WHILE READPTR <> NIL DO WITH READPTR^ DO
5700     BEGIN
5701         WRITEPS(PARSTRUCT);
5702         READPTR := NEXT
5703     END;
5704 WRITELN(DIAGFILE);
5705 WRITELN(DIAGFILE);
5706 WRITELN(DIAGFILE, 'Declarations: ');
5707 WRITELN(DIAGFILE);
5708 FOR SUBST := 1 TO MAXVARS DO IF SUBST IN DECLSUBS THEN
5709     BEGIN
5710         WRITELN(DIAGFILE, '**** R', SUBST : 1, ' ****');
5711         WRITEDECN(RDECLARATIONTABLE[SUBST]);
5712         WRITELN(DIAGFILE)
5713     END;

```

```

5714 WRITELN(DIAGFILE);
5715 WRITELN(DIAGFILE);
5716 WRITELN(DIAGFILE, 'Definitions: ');
5717 WRITELN(DIAGFILE);
5718 FOR SUBST := 1 TO MAXVARS DO IF SUBST IN DEFNSUBS THEN
5719     BEGIN
5720         WRITELN(DIAGFILE);
5721         WRITELN(DIAGFILE, '++++ R', SUBST : 1, ' +++');
5722         WRITECG(RDEFINITIONTABLE[SUBST], 0)
5723     END;
5724 WRITELN(DIAGFILE);
5725 WRITELN(DIAGFILE);
5726 WRITELN(DIAGFILE, 'Multipliers:');
5727 FOR MULT := 1 TO MAXVARS DO IF MULT IN DECLMULT THEN
5728     BEGIN
5729         WRITELN(DIAGFILE);
5730         WRITELN(DIAGFILE, ':::: M', MULT:1, '::::');
5731         WRITE(DIAGFILE, 'Values :');
5732         WRITEFREQ(MDEFINITIONTABLE[MULT]);
5733         WRITELN(DIAGFILE);
5734         WRITELN(DIAGFILE)
5735     END;
5736 WRITELN(DIAGFILE);
5737 WRITELN(DIAGFILE)
5738 END;
5739           { of Procedure OUTINTREP
5740 *****}
5741
5742
5743
5744 {-----}
5745 PROCEDURE TIDYINTREP;
5746
5747 { Deletes RDECLARATIONTABLE and RDEFINITIONTABLE, along with the latter's
5748   pendant gates. Then runs down from IRLISTTOP, deleting those PSs without
5749   parent gates, their child gates. If DIAGNOSTICS is TRUE, then a list
5750   of PSs is output, with their PSNOs, and an indication of whether or
5751   not they have been deleted. }
5752
5753 VAR DECLPTR      : PPSLIST;

```

```

5754      PMPTR      : PMDECLIST;
5755      MULT       : MULTIPLIER;
5756      SUBST      : SUBSTITUENT;
5757      OLDECTR,
5758      NUMPSS,
5759      NUMDESTROYED : INTEGER;
5760
5761
5762
5763      PROCEDURE DESTROYCG(VAR COMBINBAR : PCOMBINLIST);
5764
5765      VAR COMBINPTR : PCOMBINLIST;
5766          ALTERNPTR : PALTERNLIST;
5767
5768      BEGIN
5769      WHILE COMBINBAR <> NIL DO
5770          BEGIN
5771              COMBINPTR := COMBINBAR^.NEXT;
5772              IF COMBINBAR^.BOTTOMBAR
5773                  THEN BEGIN
5774                      DISPOSE(COMBINBAR, FALSE);
5775                      ECTRSIZE := ECTRSIZE - 24
5776                  END
5777                  ELSE BEGIN
5778                      WHILE COMBINBAR^.ALTERNATIVES <> NIL DO
5779                          BEGIN
5780                              DESTROYCG(COMBINBAR^.ALTERNATIVES^.COMBINATION);
5781                              ALTERNPTR := COMBINBAR^.ALTERNATIVES^.NEXT;
5782                              DISPOSE(COMBINBAR^.ALTERNATIVES);
5783                              ECTRSIZE := ECTRSIZE - 4;
5784                              COMBINBAR^.ALTERNATIVES := ALTERNPTR
5785                          END;
5786                              DISPOSE(COMBINBAR, TRUE);
5787                              ECTRSIZE := ECTRSIZE - 11
5788                          END;
5789                  COMBINBAR := COMBINPTR
5790          END
5791      END;
5792
5793

```

```

5794
5795 BEGIN {Body of TIDYINTREP}
5796 OLDECTR := ECTRSIZE;
5797 FOR SUBST := 1 TO MAXVARS DO IF SUBST IN DEFNSUBS THEN
5798     BEGIN
5799         REPEAT
5800             DECLPTR := RDECLARATIONTABLE[SUBST]^NEXT;
5801             DISPOSE(RDECLARATIONTABLE[SUBST]);
5802             RDECLARATIONTABLE[SUBST] := DECLPTR
5803             UNTIL RDECLARATIONTABLE[SUBST] = NIL;
5804
5805             DESTROYCG(RDEFINITIONTABLE[SUBST])
5806         END;
5807
5808     FOR MULT := 1 TO MAXVARS DO IF MULT IN DEFNMULT THEN
5809         REPEAT
5810             PMPTR := MDECLARATIONTABLE[MULT]^NEXT;
5811             DISPOSE(MDECLARATIONTABLE[MULT]);
5812             MDECLARATIONTABLE[MULT] := PMPTR
5813             UNTIL MDECLARATIONTABLE[MULT] = NIL;
5814
5815     NUMPSS := 0;
5816     NUMDESTROYED := 0;
5817     IRLISTBOT := INTERNALREP.PSLIST;
5818     WHILE IRLISTBOT <> NIL DO WITH IRLISTBOT^ DO
5819         BEGIN
5820             NUMPSS := NUMPSS + 1;
5821             IF DIAGNOSTICS
5822                 THEN WRITE(DIAGFILE, NUMPSS:6, ADDRESSOF(PARSTRUCT^): 12);
5823             IF (PARSTRUCT^.PARENTGATE = NIL) AND (PARSTRUCT <> INTERNALREP.CONSTANTPART)
5824                 THEN BEGIN
5825                 DESTROYCG(PARSTRUCT^.CHILDGATE);
5826                 CASE PARSTRUCT^.PSVARIETY OF
5827                     UNKNOWN : BEGIN
5828                         DISPOSE(PARSTRUCT, UNKNOWN);
5829                         ECTRSIZE := ECTRSIZE - 6
5830                     END;
5831                     DUMMY : BEGIN
5832                         DISPOSE(PARSTRUCT, DUMMY);
5833                         ECTRSIZE := ECTRSIZE - 8

```

```

5834          END;
5835          SPECIFIC : BEGIN
5836              DISPOSE(PARSTRUCT, SPECIFIC);
5837              ECTRSIZE := ECTRSIZE - 70
5838          END;
5839          GENERIC  : BEGIN
5840              DISPOSE(PARSTRUCT, GENERIC);
5841              ECTRSIZE := ECTRSIZE - 50
5842          END;
5843          OTHER    : BEGIN
5844              DISPOSE(PARSTRUCT, OTHER);
5845              ECTRSIZE := ECTRSIZE - 22
5846          END
5847      END;
5848      NUMDESTROYED := NUMDESTROYED + 1;
5849      IF DIAGNOSTICS THEN WRITELN(DIAGFILE, ' DESTROYED')
5850  END
5851      ELSE IF DIAGNOSTICS THEN WRITELN(DIAGFILE);
5852      IRLISTBOT := NEXT
5853  END;
5854  NUMPSS := NUMPSS - NUMDESTROYED;
5855  WRITELN('ECTR occupies ', ECTRSIZE : 5, ' words, in ', NUMPSS : 3, ' partial structures. ');
5856  WRITELN(NUMDESTROYED : 2, ' partial structures (', OLDECTR - ECTRSIZE : 5, ' words) were reclaimed. ');
5857  IF DIAGNOSTICS THEN
5858      BEGIN
5859          WRITELN(DIAGFILE);
5860          WRITELN(DIAGFILE);
5861          WRITELN(DIAGFILE);
5862          WRITELN(DIAGFILE)
5863      END
5864  END;      { of TIDYINTREP
5865  -----}
5866
5867
5868
5869
5870
5871  BEGIN          (* Body of Procedure INTERPRET *)
5872  INITIALISE;
5873  NEXTTOKEN;

```

```

5874 WHILE CHECKDELIM([GINPUT, GQUERY]) = INVALIDTOKEN DO ERROR(11,0);
5875 INTERNALREP.QUERYSTRUCTURE := TOKEN.DELIMVAL = GQUERY;
5876 NEXTTOKEN;
5877 WHILE TOKEN.NATURE <> INTEGRAL DO ERROR(16,0);
5878 INTERNALREP.REFNUMBER := TOKEN.INTEGVAL;
5879 NEXTTOKEN;
5880 WHILE CHECKDELIM([GSD]) <> GSD DO ERROR(12,0);
5881 WITH INTERNALREP DO
5882     BEGIN
5883         READSD(CONSTANTPART, INPUTMODE=TERMINAL);
5884         NEW(PSLIST);
5885         ECTRSIZE := ECTRSIZE + 4;
5886         PSLIST^.PARSTRUCT := CONSTANTPART;
5887         PSLIST^.NEXT := NIL;
5888         IRLISTBOT := PSLIST
5889     END;
5890 REPEAT
5891     STATEMENT;
5892     WHILE CHECKDELIM([GSEMI, GPERIOD]) = INVALIDTOKEN DO ERROR(24,0);
5893     IF TOKEN.DELIMVAL = GPERIOD THEN CHECKALLDONE
5894 UNTIL TOKEN.DELIMVAL = GPERIOD;
5895 RECORDMULTS;
5896 WRITELN;
5897 IF CONDITIONSPRESENT THEN WRITELN('(Conditions not yet implemented)');
5898 WRITELN;
5899 WRITELN('Generic Structure ', INTERNALREP.REFNUMBER : 6, ' accepted. ');
5900 IF DIAGNOSTICS THEN OUTINTREP;
5901 TIDYINTREP;
5902 WRITELN
5903 END;

```

## APPENDIX 4

### Line Number Index to Routines in the GENSAL Interpreter Program (Appendix 3)

PROCEDURE	ADDCOMBINPSS	: 4372
PROCEDURE	ADDCOMBSUBS	: 4928
PROCEDURE	ADDDEFNTABLE	: 4998
PROCEDURE	ADDFURTHERSUBTN	: 4571
PROCEDURE	ADDINTS	: 35
PROCEDURE	ADDITEM	: 5150
PROCEDURE	ADDPARALT	: 4343
PROCEDURE	ADDTOLIST	: 1012
PROCEDURE	ADDZERO	: 4119
FUNCTION	ALLWITHINLIMITS	: 2115
FUNCTION	ALREADYINLIST	: 4323
PROCEDURE	ALTERCONNBONDS	: 3154
PROCEDURE	ALTNTVE	: 4304
PROCEDURE	ALTNVLIST	: 2371
FUNCTION	ASSGNTOP	: 4687
PROCEDURE	ASSIGNMENTSTMNT	: 4678
FUNCTION	BONDHECK	: 2772
FUNCTION	BONDVAL	: 1456
FUNCTION	CHECK	: 468
PROCEDURE	CHECKALLDONE	: 5364
PROCEDURE	CHECKALLWITHIN	: 1234
PROCEDURE	CHECKCOMBPOSNS	: 3447
PROCEDURE	CHECKCOMPATABILITY	: 4740
FUNCTION	CHECKDELIM	: 2044
PROCEDURE	CHECKEARLIERDEFN	: 1680
PROCEDURE	CHECKINCLUDED	: 1257
FUNCTION	CHECKINT	: 581
PROCEDURE	CHECKPOSNS	: 3372
PROCEDURE	CHECKVALIDINT	: 2316
PROCEDURE	CMPDSTMNT	: 5319
FUNCTION	COMBINEDPOSITIONS	: 3411
PROCEDURE	COMBINEVALUES	: 5167
PROCEDURE	COMPARELISTS	: 1207
PROCEDURE	CONCATENATETERMS	: 2611
PROCEDURE	CONDITION	: 5269
PROCEDURE	COPYALTBAR	: 1090
PROCEDURE	COPYCOMBAR	: 1078
PROCEDURE	COPYCOMBAR	: 1278
PROCEDURE	COPYLIST	: 5118
FUNCTION	COPYLIST	: 2532
FUNCTION	COPYPS	: 1032
PROCEDURE	DECLAREMULT	: 1826
PROCEDURE	DECLARESUBST	: 1169
PROCEDURE	DECODECT	: 61
FUNCTION	DEFNTABLEENTRY	: 2700
PROCEDURE	DELETEGENSAL	: 53

PROCEDURE DESTROY	: 19
PROCEDURE DESTROYCG	: 5763
PROCEDURE DIVIDELINE	: 1985
PROCEDURE ELEMENT	: 3214
PROCEDURE ENCODECT	: 72
PROCEDURE ENTERCOMBIN	: 1371
PROCEDURE ERROR	: 776
FUNCTION EXTRACTINT	: 601
FUNCTION EXTRALAYER	: 3950
PROCEDURE FAILURE	: 317
FUNCTION FINDBOTTOM	: 4450
PROCEDURE FINDCONNECTIONS	: 3674
PROCEDURE FINDFIRST	: 2429
PROCEDURE FINDNOMEN	: 543
PROCEDURE FINDNONAPICPOSNS	: 2878
PROCEDURE FINDPOSITIONS	: 2720
PROCEDURE FINDSECOND	: 2445
PROCEDURE GETAVAILABLEPOSITIONS	: 849
PROCEDURE GETBINFO	: 4883
PROCEDURE GETCHILDPOSITIONS	: 2749
PROCEDURE GETCOMBPOSNS	: 3067
PROCEDURE GETLIMITPOSITIONS	: 4151
PROCEDURE GETMARKEDPOSNS	: 2799
PROCEDURE GETPOSNS	: 1711
PROCEDURE GETSETPOSNS	: 1726
PROCEDURE GETSPSPARAMS	: 3597
PROCEDURE GETTOKEN	: 357
PROCEDURE GOTOCOMMAND	: 13
PROCEDURE GROUPRANGE	: 2282
PROCEDURE HNUMBER	: 1577
FUNCTION HYDROGENPS	: 2830
PROCEDURE IFSTATEMENT	: 5300
FUNCTION INCREASING	: 2102
FUNCTION INDEPENDENT	: 1496
PROCEDURE INITIALISE	: 198
PROCEDURE INTEGERRANGE	: 2062
PROCEDURE INTERPRET	: 125
PROCEDURE INTSET	: 2253
PROCEDURE LISTPARAMS	: 99
PROCEDURE LISTPOSNS	: 974
FUNCTION LISTSMATCH	: 4771
FUNCTION LMAGCHECKS	: 4508
FUNCTION LMAGNOCHECKS	: 4480
PROCEDURE LOOKAHEAD	: 746
FUNCTION MAGNITUDE	: 816
FUNCTION MINBOND	: 867
FUNCTION MINPARENTBOND	: 880
PROCEDURE MODIFYCHILDPOSITIONS	: 3045
PROCEDURE MODIFYGATEPOSITIONS	: 3796
PROCEDURE MULTASSIGNMENT	: 5106
FUNCTION NEEDTOCHECK	: 3308
FUNCTION NEWCOMBAR	: 3256
FUNCTION NEWFREQ	: 210
FUNCTION NEWPARENTPSLIST	: 3561
PROCEDURE NEXTTOKEN	: 726

FUNCTION	NODENATURE	: 1434
FUNCTION	NORECORD	: 81
FUNCTION	NOVARIABLESUBTN	: 4106
FUNCTION	NUMOFCONNS	: 1659
FUNCTION	ORIGINALPOSNS	: 3338
PROCEDURE	OUTINTREP	: 5416
FUNCTION	PARAMETER	: 3652
PROCEDURE	PARAMETERLIST	: 3631
FUNCTION	POINTERLIST	: 4868
PROCEDURE	POSITIONSET	: 2416
PROCEDURE	POSNCOMBINATION	: 2463
FUNCTION	PPOSNS	: 4434
PROCEDURE	PRINTNOM	: 46
PROCEDURE	PROCESSCT	: 1415
PROCEDURE	PROGERROR	: 9
FUNCTION	PSNO	: 5426
PROCEDURE	RANGEFRAGMENT	: 2138
PROCEDURE	READCONGENERES	: 1530
PROCEDURE	READFELDMANN	: 118
PROCEDURE	READLINE	: 369
PROCEDURE	READSD	: 1948
FUNCTION	RECORDHELD	: 2656
PROCEDURE	RECORDMULTS	: 5394
PROCEDURE	REDUCE	: 2559
PROCEDURE	REDUCEECTR	: 337
PROCEDURE	REJECT	: 1512
PROCEDURE	RESTRICTSTMNT	: 5283
PROCEDURE	REVISELIMITS	: 4721
PROCEDURE	SELECTOR	: 2341
PROCEDURE	SETCOMBARS	: 3298
PROCEDURE	SETCONN BONDS	: 1120
PROCEDURE	SETINTS	: 2222
PROCEDURE	SETPARENTGATE	: 3231
FUNCTION	SPLITLINE	: 1964
FUNCTION	SPSVARIETY	: 108
PROCEDURE	STATEMENT	: 5295
PROCEDURE	STATEMENT	: 5334
PROCEDURE	SUBSTASSIGNMENT	: 5077
PROCEDURE	SUBSTASVALUE	: 4036
PROCEDURE	SUBSTCOMBINATION	: 4799
PROCEDURE	SUBSTGROUP	: 4707
FUNCTION	SUBSTNAME	: 1480
FUNCTION	SUMFILIALS	: 930
FUNCTION	TERMREAD	: 90
FUNCTION	THISWAYROUND	: 2848
PROCEDURE	TIDYINTREP	: 5745
PROCEDURE	TRACEDOWNGATE	: 3087
PROCEDURE	TRANSLATENOMEN	: 3778
PROCEDURE	UPDATEPARALTCONNS	: 2396
PROCEDURE	UPDATEPPSCONNS	: 1149
PROCEDURE	USERPARAMETER	: 3709
FUNCTION	VALIDSUBST	: 3994
FUNCTION	WITHINLIMITS	: 2076
PROCEDURE	WRITEALTERNES	: 5541
PROCEDURE	WRITECG	: 5535

APPENDIX 4:

PROGRAM INDEX

PROCEDURE WRITECG	: 5595
PROCEDURE WRITECOMBIN	: 5561
PROCEDURE WRITECONNS	: 5493
PROCEDURE WRITECT	: 5611
PROCEDURE WRITEDECN	: 5672
PROCEDURE WRITEFREQ	: 5449
PROCEDURE WRITEMESSAGE	: 281
PROCEDURE WRITEPGS	: 5515
PROCEDURE WRITEPOSNS	: 5470
PROCEDURE WRITEPS	: 5641

**APPENDIX 5**

**GLOBAL DECLARATIONS FOR GENPROG**

```

CONST MAXCT      = 32;      { Size of single connection table }
  TERMLENGTH    = 32;      { Length of nomenclatural terms }
  MAXVARS       = 63;      { Number of substituents or multipliers }
  MAXCONGENERS  = 6;       { Number of congeners }
  MAXPACKETS    = 32;
  MAXBITS       = 32;
  MAXSCREENS    = 1024;
  MAXLENGTH     = 100;     { Length of GENESIS command line, or Gensal line }
  CTFLAG        = '{';
  GENEXFLAG     = '|';
  HSTFLAG       = '}';
  CONTNFLAG     = '\';
  ENDGENFLAG    = '~';
  NOTSET        = -1;     { Indicator for INTRECORD.TOPRANGE }

```

```

TYPE TCOMMAND   =(CGENSAL, CFILE, CSAVE, CLIST, CSEARCH, CPRINT, CEDIT, CDRAW,
  CDIAGNOSE, CDICT, CNEWTERM, CSYNONYM, CSTOP, CRUN, CCURRENT,
  CFORWARD, CBACK, CTOP, CEND, CDELETE, CINSERT, CLOCATE, CHELPEDIT,
  CQUIT);
STRING14        = PACKED ARRAY[1..14] OF CHAR;
LINESTRING      = PACKED ARRAY[1..MAXLENGTH] OF CHAR;
PLINELIST       = ^LINELIST;
LINELIST        = RECORD
  LINE : LINESTRING;
  NEXT,
  LAST : PLINELIST
END;

STRING4         = PACKED ARRAY[1..4] OF CHAR;

USERTYPE        = PACKED RECORD
  FSAUTH,
  UPDAUTH,
  SWEEP : BOOLEAN;
  NAME : STRING4
END;

PDOUBLIST       = ^DOUBLIST;

```

```

DOUBLIST = RECORD
    FIRST,
    SECOND : INTEGER;
    NEXT   : PDOUBLIST
END;

INTRECORD = RECORD
    SUBRANGES : PDOUBLIST;
    TOPRANGE  : INTEGER
END;

PINTEGSET = ^INTEGSET;
INTEGSET  = SET OF 0..MAXVARS;
PTGROUPMEMS = ^TGROUPMEMS;
TGROUPMEMS = RECORD
    CASE COMBINED      : BOOLEAN OF
        TRUE  : (COMBMEMS : PDOUBLIST);
        FALSE : (MEMBERS  : INTEGSET)
    END;

BONDORDER =(NOTSPECIFIED, ANY, CHAIN, RING, SINGLE, DOUBLE, TRIPLE,
    CHAISING, CHAIDOUB, CHAITRIP, CHAITAUT, RINGSING, RINGDOUB,
    RINGTRIP, AROMATIC, RINGTAUT);

TCONNS = NOTSET..2;
TCONNBONDS = RECORD
    CASE CONNECTIONS : TCONNS OF
        NOTSET,
        0           : ();
        1           : (BOND : BONDORDER);
        2           : (BONDA,
            BONDB : BONDORDER)
    END;

TSELECTMODE =(INDEPENDENT, ALLSAME, ALLDIFF, NOTALLSAME, NOTALLDIFF);

PTRPSTYPE = ^PSTYPE;

PCOMBINLIST = ^COMBINLIST;
PALTERNLIST = ^ALTERNLIST;

COMBINLIST = RECORD
    PARENTPOSITIONS : PTGROUPMEMS;
    FREQUENCY       : INTRECORD;

```

```

NEXT          : PCOMBINLIST;
CASE BOTTOMBAR : BOOLEAN OF
  TRUE        : (CHILDPS      : PTRPSTYPE;
                 CHILDPOSITIONS : TGROUPMEMS;
                 CONNBONDS     : TCONNBONDS);
  FALSE       : (ALTERNATIVES : PALTERNLIST)
END;

```

```

ALTERNLIST = RECORD
  COMBINATION : PCOMBINLIST;
  NEXT       : PALTERNLIST
END;

```

```

PPARENTLIST = ^PARENTLIST;
PARENTLIST  = RECORD
  CHILDPOSITIONS,
  PARENTPOSITIONS : TGROUPMEMS;
  PARENTPS       : PTRPSTYPE;
  CONNBONDS     : TCONNBONDS;
  NEXT          : PPARENTLIST
END;

```

```

RELATIVES = (NONE, FRATERNAL, PARENTAL, FILIAL);
ATOMNUMBER = 0..MAXCT;
SUBSTITUENT = 0..MAXVARS;

```

```

CONGARRAY = ARRAY [1..MAXCONGENER] OF
RECORD
  BOND : BONDORDER;
  CASE RELATIONSHIP : RELATIVES OF
    NONE,
    PARENTAL,
    FILIAL : ();
    FRATERNAL : (ROWNUM : ATOMNUMBER)
END;

```

```

STRING2 = PACKED ARRAY[1..2] OF CHAR;
NUMCONGENER=0..MAXCONGENER;

```

```

ROW = RECORD

```

```
      CHARGE      : -9..9;
      HYDROGENS   : NUMCONGENERS;
      CONGENERS   : CONGARRAY;
      CASE ATOMICROW : BOOLEAN OF
        TRUE      : (ATOM      : STRING2);
        FALSE     : (NAME      : SUBSTITUENT;
                    VALUES   : PCOMBINLIST)
      END;

CTTYPE      = ARRAY[1..MAXCT] OF ^ROW;

TPARAMETERS =(ATOMCOUNT, TBRANCH, QBRANCH, EUNSATURATION, YUNSATURATION,
              RINGCOUNT, RINGATOMS, RINGSUBSTITUTION, RINGFUSIONS,
              RINGAROMATIC, HETEROATOM);

TPARAMLIST  = ARRAY[TPARAMETERS] OF INTRECORD;

TPSVARIETY  =(DUMMY, UNKNOWN, SPECIFIC, GENERIC, OTHER);
STRING32    = PACKED ARRAY[1..32] OF CHAR;

PSTYPE     = RECORD
  VISITED      : BOOLEAN;
  CHILDGATE    : PCOMBINLIST;
  PARENTGATE   : PPARENTLIST;
  CASE PSVARIETY : TPSVARIETY OF
    DUMMY      : (SUBSTNAME : SUBSTITUENT);
    UNKNOWN    : ();
    SPECIFIC   : (CT        : CTTYPE);
    GENERIC    : (PARAMLIST : TPARAMLIST);
    OTHER      : (TERM      : STRING32)
  END;

MULTIPLIER  = 0..MAXVARS;

FELDROW     = RECORD
  CHEM : STRING4 (* the atomic symbol, R group or * *);
  CHGE : -9..9;
  MULT : MULTIPLIER;
  AR   : ARRAY [1..MAXCONGENERS] OF ATOMNUMBER (* the congeners *)
  END;
```

```

BONDROW      = RECORD
              NODE1,
              NODE2 : ATOMNUMBER;
              BOND  : 1..16
            END;
TFELDMODE    =(NEWDIAGRAM, OLDDIAGRAM, NUMBERDRAW, NUMBERLESSDRAW);

PIRLIST      = ^TIRLIST;
TIRLIST      = RECORD
              PARSTRUCT : PTRPSTYPE;
              NEXT      : PIRLIST
            END;

```

```

TINTERNALREP= RECORD
              REFNUMBER      : INTEGER;
              QUERYSTRUCTURE : BOOLEAN;
              CONSTANTPART   : PTRPSTYPE;
              PSLIST         : PIRLIST;
            END;

```

```

VAR DIAGFILE      : TEXT;           { Diagnostics file variable }
    TOPOGMFILE    : TEXT;           { Grammar file variable }
    USERFILE      : FILE OF USERTYPE;
    USER          : USERTYPE;
    DIAGFIL       : ALFA;           { Diagnostics file name }
    DIAGNOSTICS,  : ALFA;           { Diagnostics file indicator }
    STRUCTURECOMPLETED : BOOLEAN;

BUFFER           : LINestring;      { GENESIS or EDITOR command line or Gensal line }
N                : 0..MAXLENGTH;   { Character counter for BUFFER }
WORKSPACE        : PLINELIST;       { Static pointer to held GENSAL }
INTERNALREP      : TINTERNALREP;
INSERTGENEX      : PLINELIST;       { Lines of SPSfile Gensal expression }
SPSPARAMLIST     : TPARAMLIST;      { SPSfile parameter list }
FELDCT           : ARRAY [1..MAXCT] OF FELDROW; { The Feldmann connection table }
FELDBD           : ARRAY [1..MAXCT] OF BONDROW; { The Feldmann bonding table }
FELDMODE         : TFELDMODE;       { Calling mode for FELDMN }
NUMOFNODES,     : INTEGER;          { Number of nodes in the Feldmann connection table }

```

```
NUMOFBONDS      : ATOMNUMBER;           { Number of bonds in the Feldmann bonding table }
FELDFIL         : ALFA;                 { Feldmann transfer file }
```

```
{-----}
PRIME APPLICATIONS LIBRARY ROUTINES }
```

```
FUNCTION CLOS$(UNIT : SHORTINT): BOOLEAN;
EXTERN;
```

```
FUNCTION OPN$(OPNKEY : SHORTINT;
              NAME   : STRING14;
              NAMLEN,
              UNIT,
              VERKEY,
              WTIME,
              RETRYS : SHORTINT): BOOLEAN;
EXTERN;
```

```
FUNCTION POSN$(POSKEY,           { 1=A$ABS }
              UNIT : SHORTINT;
              POS  : INTEGER) : BOOLEAN;
EXTERN;
```

```
{-----}
EXTERNAL FORTRAN SUBROUTINES
  Loaded in INOUTSUBS}
```

```

PROCEDURE GETLIN(VAR LINE : LINESTRING);
EXTERN;
{ FORTRAN subroutine to obtain a single line from the file already open on
  unit 1 and positioned at the correct place. }

```

```

PROCEDURE FELDMN(VAR FELDMODE : TFELDMODE;
                 VAR FELDFIL  : ALFA);
EXTERN;
{ Displays a structure diagram, for which the connection
  table is in the file FELDFIL }

```

```

PROCEDURE GETNOM(VAR TERM : STRING32;
                 VAR ADDR : INTEGER);
EXTERN;
{ FORTRAN subroutine to obtain the next TERM and ADDR from the file
  SPSDICT, which is already open on unit 1. }

```

```
{*****}
```

```

PROCEDURE ADDINTS (VAR PTR1      : PDOUBLIST;
                  LOWER, UPPER : INTEGER);

VAR PTR2 : PDOUBLIST;

BEGIN
  IF PTR1 = NIL
    THEN PTR2 := NIL
    ELSE WITH PTR1^ DO
      IF (SECOND = LOWER-1) OR (SECOND = LOWER)
        THEN SECOND := UPPER
        ELSE BEGIN

```

```

                PTR2 := PTR1;
                PTR1 := NIL
            END;
IF PTR1 = NIL
    THEN BEGIN
        NEW(PTR1);
        WITH PTR1^ DO
            BEGIN
                FIRST := LOWER;
                SECOND := UPPER;
                NEXT := PTR2
            END
        END
    END;

```

```

PROCEDURE PRINTNOM(NOMENVAL : STRING32);

```

```

    VAR M : 1..32;

```

```

    BEGIN

```

```

        FOR M := 1 TO 32 DO

```

```

            IF NOMENVAL[M] <> ' ' THEN WRITE(NOMENVAL[M])
        END;

```

```

PROCEDURE DELETEDGENSAL(VAR LINE1 : PLINELIST);

```

```

    VAR LINE2 : PLINELIST;

```

```

    BEGIN

```

```

        WHILE LINE1 <> NIL DO

```

```

            BEGIN

```

```

                LINE2 := LINE1^.NEXT;

```

```

                DISPOSE(LINE1);

```

```

                LINE1 := LINE2
            END
        END
    END;

```

{-----}

```
PROCEDURE DECODECT (VAR CTLINE : PLINELIST;  
                    DISPLAYING : BOOLEAN);
```

```
VAR CHPOSN      : 0..MAXLENGTH;    { Character position in LINE }  
    NODE        : ATOMNUMBER;      { Loop counter }  
    M,          :                   { Miscellaneous counter }  
    SPACE       : INTEGER;         { Ordinal value offset }
```

```
FUNCTION NEXTCH : CHAR;
```

```
{ Returns the next character in the string, taking new lines when necessary }
```

```
BEGIN
```

```
IF CHPOSN=MAXLENGTH
```

```
THEN BEGIN
```

```
    CTLINE := CTLINE^.NEXT;
```

```
    CHPOSN := 2 { First character in each line is omitted (CONTNFLAG) }
```

```
END
```

```
ELSE CHPOSN := CHPOSN + 1;
```

```
NEXTCH := CTLINE^.LINE [CHPOSN]
```

```
END;
```

```
BEGIN { Body of DECODECT }
```

```
SPACE := ORD(' ');
```

```
CHPOSN := 1; { first character in string is omitted }
```

```
NUMOFNODES := ORD(NEXTCH) - SPACE;
```

```
FOR NODE := 1 TO NUMOFNODES DO WITH FELDCT[NODE] DO
```

```
    BEGIN
```

```
        FOR M := 1 TO 4 DO CHEM[M] := NEXTCH;
```

```
        CHGE := ORD(NEXTCH) - SPACE - 9;
```

```
        MULT := ORD(NEXTCH) - SPACE;
```

```
        FOR M := 1 TO MAXCONGENERS DO AR[M] := ORD(NEXTCH) - SPACE
```

```

END;
NUMOFBONDS := ORD(NEXTCH) - SPACE;
FOR M := 1 TO NUMOFBONDS DO WITH FELDBD[M] DO
BEGIN
  NODE1 := ORD(NEXTCH) - SPACE;
  NODE2 := ORD(NEXTCH) - SPACE;
  BOND  := ORD(NEXTCH) - SPACE
END;
IF DISPLAYING THEN
BEGIN
  REWRITE(OUTPUT,FELDFIL);
  WRITELN(NUMOFNODES : 3);
  FOR NODE := 1 TO NUMOFNODES DO WITH FELDCT[NODE] DO
  BEGIN
    WRITE(CHEM, CHGE:2, ' ');
    IF MULT=0 THEN WRITE(' ')
      ELSE WRITE('M',MULT:3);
    FOR M := 1 TO MAXCONGENERS DO IF AR[M] <> 0 THEN WRITE(AR[M] : 3);
    WRITELN
  END;
  WRITELN(NUMOFBONDS : 3);
  FOR M := 1 TO NUMOFBONDS DO WITH FELDBD[M] DO
    WRITELN(NODE1 : 3, NODE2 : 3, BOND : 3);
  REWRITE(OUTPUT,'@TTY');
  FELDMODE := NUMBERDRAW;
  FELDMN(FELDMODE,FELDFIL)
END
END;
{-----}

```

```

{-----}

```

```

PROCEDURE ENCODECT(VAR CTLINE : PLINELIST);

```

```

VAR CHPOSN : 0..MAXLENGTH; { Charcter position in line }
    NODE    : ATOMNUMBER;   { Loop counter }
    M,      :                 { Miscellaneous counter }
    SPACE   : INTEGER;      { Ordinal value offset }

```

```

PROCEDURE STORECHAR (CH : CHAR);

{ Stores CH in the next position in the character string, taking new lines when
  necessary }

VAR NEWLINE : PLINELIST;

BEGIN
  IF CHPOSN=MAXLENGTH
    THEN BEGIN
      NEW(NEWLINE);
      NEWLINE^.NEXT := NIL;
      NEWLINE^.LAST := CTLINE;
      CTLINE^.NEXT := NEWLINE;
      CTLINE := NEWLINE;
      CTLINE^.LINE[1] := CONTNFLAG;
      CHPOSN := 2
    END
  ELSE CHPOSN := CHPOSN+1;
  CTLINE^.LINE [CHPOSN] := CH
END;

```

```

BEGIN {Body of ENCODECT}
SPACE := ORD(' ');
CHPOSN := 0;
STORECHAR(CTFLAG); { Connection table indicator flag }
STORECHAR(CHR(NUMOFNODES + SPACE));
FOR NODE := 1 TO NUMOFNODES DO WITH FELDCT[NODE] DO
  BEGIN
    FOR M := 1 TO 4 DO STORECHAR(CHEM[M]);
    STORECHAR(CHR(CHGE+9+SPACE));
    STORECHAR(CHR(MULT+SPACE));
    FOR M := 1 TO MAXCONGENERS DO STORECHAR(CHR(AR[M]+SPACE))
  END;
STORECHAR(CHR(NUMOFBONDS + SPACE));
FOR M := 1 TO NUMOFBONDS DO WITH FELDBD[M] DO

```

```

BEGIN
  STORECHAR (CHR (NODE1+SPACE));
  STORECHAR (CHR (NODE2+SPACE));
  STORECHAR (CHR (BOND+SPACE))
END;
WHILE CHPOSN < MAXLENGTH DO STORECHAR(' ')
END;
{-----}

{-----}
PROCEDURE READSPSPARAMS (SPSTRING : PLINELIST);

VAR CH      : CHAR;
    PTR     : PDOUBLIST;
    CHPOSN  : 0..MAXLENGTH;
    PARAMETER : TPARAMETERS;

FUNCTION NEXTCH : CHAR;
{ Returns the next character in the string, taking new lines when necessary }

BEGIN
  IF CHPOSN=MAXLENGTH
    THEN BEGIN
      SPSTRING := SPSTRING^.NEXT;
      CHPOSN := 1
    END
    ELSE CHPOSN := CHPOSN + 1;
  NEXTCH := SPSTRING^.LINE [CHPOSN]
END;

BEGIN {Body of READSPSPARAMS}
  CHPOSN := 1; {first character in string is omitted (HSTFLAG) }
  FOR PARAMETER:= ATOMCOUNT TO HETEROATOM DO WITH SPSPARAMLIST[PARAMETER] DO

```

```

BEGIN
  SUBRANGES := NIL;
  TOPRANGE := ORD(NEXTCH) - ORD('0');
  CH := NEXTCH;
  WHILE CH <> ' ' DO
    BEGIN
      NEW(PTR);
      PTR^.NEXT := SUBRANGES;
      PTR^.FIRST := ORD(CH) - ORD('0');
      PTR^.SECOND := ORD(NEXTCH) - ORD('0');
      SUBRANGES := PTR;
      CH := NEXTCH
    END
  END
END;
{-----}

```

```

FUNCTION NORECORD(NOMEN      : STRING32;
                  VAR ADDRESS : INTEGER): BOOLEAN;

```

```

VAR SPSNOM : STRING32;

```

```

BEGIN
  IF NOT OPNV$(SHORT(1), 'LI2GEN>SPSDICT', SHORT(14), SHORT(1), SHORT(1), SHORT(1), SHORT(100))
    THEN PROGERROR(101); {File error - cannot open SPSDICT}
  REPEAT GETNOM(SPSNOM, ADDRESS)
  UNTIL (SPSNOM=NOMEN) OR (SPSNOM[1]=' ');
  IF NOT CLOS$(SHORT(1)) THEN PROGERROR(102); {cannot close SPSDICT}
  NORECORD := SPSNOM[1]=' '
END;

```

```

FUNCTION TERMREAD(VAR TERM : STRING32) : BOOLEAN;

```

```

VAR M, MM : 0..TERMLENGTH;

```

```

BEGIN

```

```

READLN(TERM : M);
FOR MM := 1 TO M DO IF TERM[MM] IN ['a'..'z']
  THEN TERM[MM] := CHR(ORD(TERM[MM]) - ORD('a') + ORD('A'));
TERMREAD := M>0
END;

```

```

PROCEDURE LISTPARAMS(VAR OUTFILE : TEXT;
                    PARAMLIST : TPARAMLIST);

```

```

VAR PARAMETER : TPARAMETERS;
    PTR       : PDOUBLIST;

```

```

BEGIN
FOR PARAMETER := ATOMCOUNT TO HETEROATOM DO WITH PARAMLIST[PARAMETER] DO
  IF TOPRANGE <> 0 THEN
    BEGIN
      CASE PARAMETER OF
        ATOMCOUNT      : WRITE(OUTFILE, 'C');
        TBRANCH         : WRITE(OUTFILE, 'T');
        QBRANCH         : WRITE(OUTFILE, 'Q');
        EUNSATURATION   : WRITE(OUTFILE, 'E');
        YUNSATURATION   : WRITE(OUTFILE, 'Y');
        RINGCOUNT      : WRITE(OUTFILE, 'RC');
        RINGATOMS       : WRITE(OUTFILE, 'RN');
        RINGSUBSTITUTION : WRITE(OUTFILE, 'RS');
        RINGFUSIONS     : WRITE(OUTFILE, 'RF');
        RINGAROMATIC    : WRITE(OUTFILE, 'RA');
        HETEROATOM      : WRITE(OUTFILE, 'Z');
      END;
      PTR := SUBRANGES;
      WRITE(OUTFILE, '<');
      WHILE PTR <> NIL DO WITH PTR^ DO
        BEGIN
          WRITE(OUTFILE, FIRST : 1);
          IF FIRST <> SECOND THEN WRITE(OUTFILE, '-', SECOND : 1);
          PTR := NEXT;
          IF (PTR <> NIL) OR (TOPRANGE <> NOTSET) THEN WRITE(OUTFILE, ',');
        END;
      END;
    END;
  END;
END;

```

```

        IF TOPRANGE <> NOTSET
        THEN WRITE(OUTFILE, TOPRANGE:1, '-> ')
        ELSE WRITE(OUTFILE, '> ')
    END;
WRITELN(OUTFILE);
WRITELN(OUTFILE)
END;

```

```

FUNCTION SPSVARIETY(ADDRESS      : INTEGER;
                   DISPLAYING : BOOLEAN) : TPSVARIETY;

```

```

{ Returns the variety of partial structure, whose record begins at ADDRESS in
  SPSFILE. The lines of the record in SPSFILE are in reverse order, and as they
  are read into a linked list of lines, the order is automatically put right.
  The first character of the first (in correct order) line indicates the nature
  of the partial structure. DECODECT is called to deal with connection tables
  (with DISPLAYING as its parameter); homologous series terms are handled by
  READSPSPARAMS, and listed by LISPARAMS if DISPLAYING is TRUE; Gensal
  expressions are stored in INSETGENEX, and listed by LISTGENEX if DISPLAYING
  is TRUE. }

```

```

VAR SPSTRING : PLINELIST; { Lines of partial structure record }
    PARAMETER : TPARAMETERS;

```

```

BEGIN
IF NOT (OPNV$(SHORT(1), 'LI2GEN>SPSFILE', SHORT(14), SHORT(1), SHORT(1), SHORT(1), SHORT(100))
      AND POSN$(SHORT(1), SHORT(1), ADDRESS)) THEN
    PROGERROR(103); {File error - opening/positioning SPSFILE}
NEW(SPSTRING);
SPSTRING^.LAST := NIL;
SPSTRING^.NEXT := NIL;
GETLIN(SPSTRING^.LINE);
WHILE NOT (SPSTRING^.LINE[1] IN [CTFLAG, HSTFLAG, GENEXFLAG]) DO
    BEGIN
        NEW(SPSTRING^.LAST);
        SPSTRING^.LAST^.NEXT := SPSTRING;
        SPSTRING := SPSTRING^.LAST;
        SPSTRING^.LAST := NIL;
    END

```

```

    GETLIN(SPSTRING^.LINE)
  END;
IF NOT CLOS$(SHORT(1)) THEN PROGERR(104); {File error (SPSVARIETY) - closing SPSFILE}
CASE SPSTRING^.LINE[1] OF
  CTFLAG:
    BEGIN
      SPSVARIETY := SPECIFIC;
      DECODECT(SPSTRING, DISPLAYING);
      DELETEGENSAL(SPSTRING)
    END;
  GENEXFLAG:
    BEGIN
      SPSVARIETY := OTHER;
      INSERTGENEX := SPSTRING^.NEXT;
      IF DISPLAYING THEN LISTGENEX(SPSTRING^.NEXT)
    END;
  HSTFLAG:
    BEGIN
      SPSVARIETY := GENERIC;
      READSPSPARAMS(SPSTRING);
      DELETEGENSAL(SPSTRING);
      IF DISPLAYING
        THEN BEGIN
          LISTPARAMS(OUTPUT, SPSPARAMLIST);
          FOR PARAMETER := ATOMCOUNT TO HETEROATOM DO
            DESTROY(SPSPARAMLIST[PARAMETER].SUBRANGES)
        END
    END
END { of case }
END;

```

```
PROCEDURE READFELDMANN;
```

```
{ Reads the Feldmann table from FELDFIL. }
```

```

VAR CH   : CHAR;
    NODE : ATOMNUMBER;
    M    : INTEGER;

```

```
BEGIN
RESET(INPUT,FELDFIL);
READLN(NUMOFNODES);
FOR NODE := 1 TO NUMOFNODES DO WITH FELDCT[NODE] DO
  BEGIN
    FOR M := 1 TO 4 DO READ(CHEM[M]);
    READ(CHGE,CH,CH,CH);
    IF CH = 'M' THEN READ(MULT)
      ELSE MULT := 0;
    FOR M := 1 TO MAXCONGENERS DO
      IF EOLN(INPUT) THEN AR[M] := 0
        ELSE READ(AR[M]);
    READLN
  END;

READLN(NUMOFBONDS);
FOR M := 1 TO NUMOFBONDS DO WITH FELDBD[M] DO
  READLN(NODE1,NODE2,BOND);
RESET(INPUT,'@TTY')
END;
```

APPENDIX 5:

GLOBAL DECLARATIONS

## APPENDIX 6

### SAMPLE INTERPRETER SESSION

In this sample interpreter session the structure shown in Figure 3.3 is entered, with various errors being indicated by the program, and corrected by the user. After a "failure", a session using the editor corrects an erroneous structure diagram, and the whole structure is then reprocessed in non-interactive mode, before the user continues to input GENSAL statements.

Enter Command : GENSAL

```
1      18  GENSAL: INPUT 4163058
2      18  GENSAL: SD
```

FELDMANN graphics system for structure diagram input and display:

```
#
RING 5
```

```
#
ABRAN 1 1 2 1 3 1 4 1
```

```
#
SATOM 1 3
ATOM TYPE =
N
```

```
#
SATOM 6 8
ATOM TYPE =
R5
```

```
#
SATOM 7 9
ATOM TYPE =
0
```

#  
SBOND 2 7 4 9  
BOND TYPE=  
CD

#  
RING 6

#  
RING 6

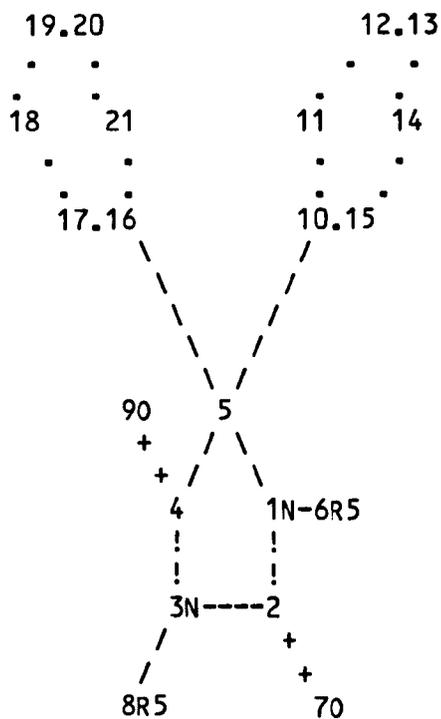
#  
ABOND 5 10

#  
ABOND 5 16

#  
ALTBD 10 11

#  
ALTBD 16 17

#  
D



#  
END  
3 702 GENSAL: R = H / SD

\*\*\*\* ERROR 23  
Integer expected.

Remainder of input line ignored

4 702 GENSAL: 5 = H / SD

FELDMANN graphics system for structure diagram input and display:

```
#
CHAIN 3

#
ABRAN 1 1 1 1

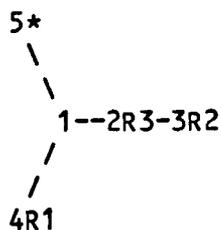
#
SATOM 2
ATOM TYPE =
R3

#
SATOM 3
ATOM TYPE =
R2

#
SATOM 4
ATOM TYPE =
R1

#
SATOM 5
ATOM TYPE =
*

#
D
```



```
#
END
5 1299 GENSAL: R1 = H / alkyl <1-7> ;
```

\*\*\*\* ERROR 24  
Unexpected symbol.

Remainder of input line ignored

```
6 1299 GENSAL: ;
7 1299 GENSAL: R1 = H / alkyl <1-7> ;
```

8 1648 GENSAL: R2 = SD

FELDMANN graphics system for structure diagram input and display:

```
#
CHAIN 2

#
ABRAN 1 1 1 1

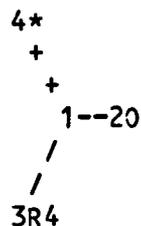
#
SATOM 2
ATOM TYPE =
0

#
SATOM 3
ATOM TYPE =
R4

#
SATOM 4
ATOM TYPE =
*

#
SBOND 1 4
BOND TYPE=
CD
```

```
#
D
```



```
#
END
```

```
**** FAILURE 42
Bond types CS and CD are incompatible.
```

Edit existing GENSAL or start again!

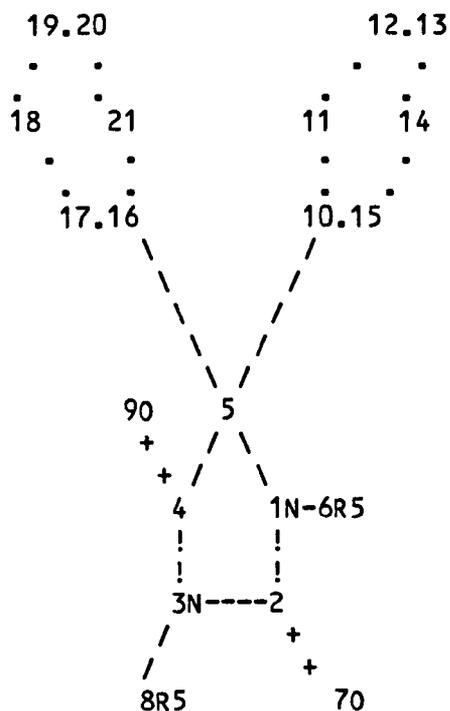
Enter Command : EDIT

[The editor session is not shown here. It involves the replacement of the erroneous double bond in the last

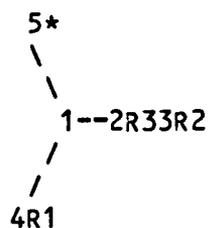
structure diagram by a single bond]

> RUN

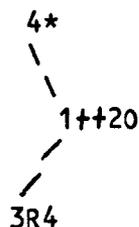
1 18 GENSAL: INPUT 4163058  
2 18 GENSAL: SD



3 702 GENSAL: R  
4 702 GENSAL: 5 = H / SD



5 1299 GENSAL:  
6 1299 GENSAL: ;  
7 1299 GENSAL: R1 = H / alkyl <1-7> ;  
8 1648 GENSAL: R2 = SD



End of stored GENSAL.  
Input at the terminal:

```
9 1887 GENSAL: ;  
10 1887 GENSAL: R3 = 0 / S;  
11 2292 GENSAL: RESTRICT <1> R5 <> H.
```

The following substituents remain undefined:  
R4

\*\*\*\* ERROR13  
";" expected.

Remainder of input line ignored

```
12 2292 GENSAL: ;  
13 2292 GENSAL: R4 = 'acyl residue of naturally-  
14 2325 GENSAL:      occurring protein amino acid'.
```

(Conditions not yet implemented)

Generic Structure 4163058 accepted.  
ECTR occupies 2124 words, in 9 partial structures.  
0 partial structures ( 279 words) were reclaimed.

Used 5.734 seconds.

Fragment generation begins.

## APPENDIX 7

### INTERPRETER ERROR MESSAGES

In these error messages the symbol # is replaced by an integer, and the symbol \$ by a character (normally part of a bond type abbreviation).

- 1) Substituent R# has not been declared.
- 2) Multiplier M# has not been declared.
- 3) Position # which is implicit in nomenclatural term is not available.
- 4) No positions available for the further substitution implicit in nomenclatural term.
- 5) # is not a valid value for this parameter.
- 6) Position # is not available for attachment in all child structures.
- 7) Substitution is not possible in position #.
- 8) Insufficient substitutable positions for # substitutions.
- 9) Not enough substituent declarations for # selective definitions.
- 10) # is too big a value for a multiplier.
- 11) "INPUT" or "QUERY" expected.

- 12) "SD" expected.
- 13) ";" expected.
- 14) ")" expected.
- 15) ";" or "END" expected.
- 16) Patent number expected.
- 17) "THEN" expected.
- 18) Substituent definition element expected.
- 19) Statement expected.
- 20) Substituent or multiplier group expected.
- 21) "<" expected.
- 22) ">" expected.
- 23) Integer expected.
- 24) Unexpected symbol.
- 25) Nomenclatural term or "" expected.
- 26) Substituent expected.
- 27) Integer range must have increasing values left to right.
- 28) Substituent values must be in the range 1 to 63.
- 29) "" expected.
- 30) Combination of doubly-connected substituents not permitted.
- 31) Connectivity incompatible with substituent(s) being defined.
- 32) Assignment operator expected.
- 33) "/" expected.
- 34) Position combination not permitted for singly-connected substitution.
- 35) No appropriate declaration for R#.
- 36) "+" expected.
- 37) Substituents in combination not declared in same partial structures.

- 38)
- 39) Positions specified previously are not available.
- 40) Error in stored GENSAL.
- 41) Structure Diagram rejected.
- 42) Bond types \$\$ and \$\$ are incompatible.
- 43) No available position in child structure for \$\$ bond.
- 44) Doubly-connected value for singly-connected substituent.
- 45) Bond size cannot be accomodated in parent structure at specified position #.
- 46) Bond size cannot be accomodated in parent structure at any position.
- 47) No positions available in parent structure for combined substitution.
- 48) Bond sizes cannot be accomodated at any pair of positions in parent structure.
- 49) Positions specified for further substitution on R# are not available.
- 50) No available positions for further substitution on R#.
- 51) Incompatible connectivities for substituents in this group.
- 52) No available positions in child structure for \$\$ and \$\$ bonds.
- 53) Number of bonds on substituent at node # does not agree with previous declaration.
- 54) More than two bonds on substituent at node #.
- 55) Illegal valency on atom at node #.
- 56) Variable-position label applied to atom at node #.
- 57) Multiplier applied to node #, which is not a substituent.

- 58) Multivalent label or hydrogen at node #.
- 59) Maximum of two apical bond labels exceeded.
- 60) Illegal pattern of AROMATIC bonds at node #.
- 61) Illegal pattern of TAUTOMERIC bonds at node #.

## BIBLIOGRAPHY

---

"When you steal from one author, it's plagiarism  
if you steal from many, it's research"

William Mizner (1876-1933)

1. Anon., "No growth in chemical literature in 1981", CAS Report (12) 10 (1982)
2. Rowlett, R.J., "Gleaning patents with Chemical Abstracts", ChemTech 9 (6) 348-9 (1979)
3. Oppenheim, C., "The patents coverage of Chemical Abstracts", Information Scientist 8 (3) 133-8 (1974)
4. Oppenheim, C., "Recent changes in patent law and their implications for information services and information scientists", Journal of Documentation 34 (3) 217-229 (1978)
5. Dodd, V.S., "Developments in patent documentation", Aslib Proceedings 31 (4) 180-190 (1979)
6. Kaback, S.M., "Chemical structure searching in Derwent's World Patent Index", Journal of Chemical Information and Computer Sciences 20 (1) 1-6 (1980)
7. Markush, E.A., "Pyrazolone dyes", U.S. Patent No. 1506316, Aug 26th 1924
8. Rosa, M.C., "Outline of practice relative to 'Markush' claims", Journal of the Patent Office Society 34 (5) 324-345 (1952)
9. Coulter, R.I., "Comments on 'alternativeness' in claims (The Rule Against Or)", Journal of the Patent Office Society 33 (11) 819-831 (1951)
10. Silk, J.A., Personal Communication to M.F. Lynch, Letter dated 19 October 1979
11. Bouman, H., "Too prolific Markush", Journal of Documentation 26 (2) 161-163 (1970)
12. Beton, J.L., "The nature and grant of chemical patents", Chemistry and Industry 298-301 (5 April 1975)

## BIBLIOGRAPHY

13. Valance, E.H., "Understanding the Markush claim in chemical patents", Journal of Chemical Documentation 1 (2) 87-92 (1961)
14. Sneed, H.M.S., Turnipseed, J.H., Turpin, R.A., "A line-formula notation system for Markush structures", Journal of Chemical Documentation 8 (3) 173-178 (1968)
15. Krishnamurthy, E.V., Lynch, M.F., "Analysis and coding of generic chemical formulae in chemical patents", Journal of Information Science 3 75-79 (1981)
16. Krishnamurthy, E.V., Lynch, M.F., "Formal description, coding and computer handling of generic formulae in chemical patents. Report on tenure of a British Library Senior Visiting Fellowship by Prof. E.V. Krishnamurthy at the Postgraduate School of Librarianship and Information Studies, University of Sheffield, February 1979." BLR&D Report No 5490, British Library (1979)
17. Geivandov, E.A., ["Language for notation of generalised structures of organic compounds containing alternative delocalised fragments (Markush Structures)"], Nauchno-tehnicheskaya Informatsiya Seriya 2 (10) 21-24, 46 (1972) [IN RUSSIAN]
18. Warr, W.A., "Software for storage and retrieval of chemical structures", presented at CNA(UK) Seminar on the Future of Chemical Documentation, Exeter, September 1982.
19. Eakin, D.R., "The ICI CROSSBOW system", in "Chemical Information Systems", eds. Ash, J.E., Hyde, E. Ellis Horwood (1975)
20. Smith, E.G., Baker, P.A., "The Wiswesser Line-Formula Chemical Notation", 3rd edition, Chemical Information Management Inc. (1972)
21. Polton, D., "Installation and experience with MACCS (Molecular Access System)", Online Review, 6 (3) 235-242 (1982)
22. Farmer, N.A., O'Hara, M.P., "CAS ONLINE: a new source of substance information from Chemical Abstracts Service", Database 3 10-25 (1980)
23. Dubois, J.E., "French national policy for chemical information and the DARC system as a potential tool of this policy", Journal of Chemical Documentation, 13 (1) 8-13 (1973)
24. Dubois, J.E., "The DARC system in chemistry", in "Computer Representation and Manipulation of Chemical Information", eds. Wipke, W.T., Heller, S.R., Feldmann,

## BIBLIOGRAPHY

- R.J., Hyde, E. Wiley (1974)
25. Dubois, J.E., "Ordered chromatic graph and limited environment concept", in "Chemical Applications of Graph Theory", ed. Balaban, A.T. Academic Press (1976)
  26. Jackson, F.T., "Markush Structures", Proceedings of the CNA(UK) Seminar on Integrated Databases for Chemical Systems, 9-11 April 1979, University of Kent at Canterbury, pp. 134-157. CNA(UK) (1981).
  27. Japan Patent Association, "Patent Information Study Team in Europe and America. Report, 31 August 1975", Derwent Publications (1976)
  28. Smith, R.G., Anderson, L.P., Jackson, S.K., "Online retrieval of chemical patent information. An overview and brief comparison of three major files", Journal of Chemical Information and Computer Sciences **17** (3) 148-157 (1977)
  29. Silk, J.A., "Present and future prospects for structural searching of the journal and patent literature", Journal of Chemical Information and Computer Sciences **19** (4) 195-198 (1979)
  30. Oppenheim, C., "Patent Information Online - a Review", Proceedings of the 5th International Online Information Meet Cunard Hotel, London, December 1981, pp. 91-99. Learned Information, Oxford (1981).
  31. Rowland, J.F.B., "Information Transfer and Use in Chemistry. Final Report of the Chemical Information Review Committee", British Library R&D Department Report 5385 (1978)
  32. Rowland, F., "How will chemical information develop?", Chemistry in Britain **14** (7) 342-4 (1978)
  33. Kaback, S.M., "A user's experience with the Derwent patent files", Journal of Chemical Information and Computer Sciences **17** (3) 143-148 (1977)
  34. Kaback, S.M., "Derwent Search Aids", Database, **5** (3) 19-21 (1982)
  35. Hyams, M., "Chemical patents information", Chemistry in Britain **6** 416-420 (1968)
  36. Deforeit, H., Caric, A., Combe, H., Leveque, S., Malka, A., Valls, J., "CORA - a semiautomatic coding system. Application to the coding of Markush formulas", Journal of Chemical Documentation **12** (4) 230-232 (1972)

## BIBLIOGRAPHY

37. Ramussen, L.E., Van Oot, J.G., "Operation of Du Pont's central patent index", Journal of Chemical Documentation 9 (4) 201-206 (1969)
38. Balent, M.Z., Emberger, J.M., "A unique chemical fragmentation system for indexing patent literature", Journal of Chemical Information and Computer Sciences 15 (2) 100-104 (1975)
39. Balent, M.Z., Lotz, J.W., "Polymers and patents don't mix - easily", Journal of Chemical Information and Computer Sciences 19 (2) 80-83 (1979)
40. Cattley, J.M., Rief, T.A., Moore, J.E., Banks, D.G., O'Leary, P.T., "Retrieving patents by weighted term searching", Chemical Engineering Progress 62 (10) 91-96 (1966)
41. Rössler, S., Kolb, A., "The GREMAS system, an integral part of the IDC system for chemical documentation", Journal of Chemical Documentation 10 (2) 128-134 (1970)
42. Internationale Dokumentationsgesellschaft für Chemie m.b.H. Fachinformationszentrum Chemie, "IDC System Overview. 2. IDC and its methods of operation", IDC (Undated)
43. Fugmann, R., "The IDC System", in "Chemical Information Systems", eds. Ash, J.E., Hyde, E., pp. 195-226, Ellis Horwood (1975)
44. Pätzcher, G., "Dokumentation vom Markush-Strukturen in Organischen Patenten" ["Documentation of Markush structures in organic patents"], in "Kleincomputer in Information und Dokumentation", Deutscher Dokumentartag 1981, Mainz, 5-8 Oktober 1981, ed. H. Strohl-Goebel, K.G. Saur, Munich (1982)
45. Mullen, A., "Informationswesen in der Chemie. Teil II. Neue Tendenzen" [Means of information in chemistry. Part II. New trends], Praxis der Naturwissenschaften Chemie 29 (8) 243-247 (1980) [IN GERMAN]
46. Oppenheim, C., "The performance of the Chemical Abstracts subject and formula indexes in retrieving compounds disclosed in chemical patent Information Scientist 9 (3) 107-111 (1975)
47. Kaindl, H., "The function of the internal database - burial ground or intelligence service", Presented at CNA(UK) Seminar on "The Future of Chemical Documentation", University of Exeter, September 1982
48. Schwartz, J.H., "Some observations concerning Chemical Abstracts' formula indexes", Journal of Chemical

Documentation 9 (3) 169-171 (1969)

49. Silk, J.A., Howarth, K.E., Waterman, J.R., Wilkins, M., "Substructure searching of the CA Registry file. Some implications of a publicly-available system", ICI Chemical Structures Task Force Private Circulation Document. Unpublished (1980)
50. Bois, R., Chaumier, J., "A comparative analysis of the DARC system and the information and documentation system of the IDC", World Patent Information 2 (2) 61-66 (1980)
51. Dubois, J.E., Veillard, H., Panaye, A., "Systeme DARC. XV. Theorie de generation-description VI. Les structures ambigues: description et chainage par un graphe fictif documentaire", Bulletin de la Societe Chimique de France. Part 2 1996-2002 (1973)
52. Berte, M., Delaet, F., "Recherche documentaires automatiques dans le domaine des Prostaglandines a l'aide du Systeme DARC" [Automatic document search in the Prostaglandin field by means of the DARC system], Bulletin des Societes Chimiques Belges 88 (3) 175-191 (1979). [In French, English Abstract].
53. Gay, J.-P., "The DARC system", Presented at CNA(UK) Seminar on the "Future of Chemical Documentation", University of Exeter, September 1982
54. Dyson, G.M., "Generic (or Markush) groups in notation and search programs, with particular reference to patents", Information Storage and Retrieval 2 59-71 (1964).
55. Hayward, H.W., Tauber, S.J., "The HAYSTAQ experiment", in Proceedings of the Fifth Annual Meeting of the Commit for International Cooperation in Information Retrieval Among Examining Patent Offices (ICIREPAT), held London, 31 August 10 September 1965, pp 337-50. Published Thompson, Washington D.C. (1966)
56. Tauber, S.J., "Digital handling of chemical structures and associated information", in Proceedings of 20th National Conference of the Association for Computing Machinery, Ohio, August 1965, pp. 206-16 ACM Publication P-65, Lewis Winner, New York (1965)
57. Kirby, C.E.L., Anderson, R.K., Tauber, S.J., "Tentative connection table word format for Markush-type generic chemical structures", Appendix C (Working Paper 66-3 (27 January 1966)) to "Progress in techniques for manipulating and organising chemical information", National Bureau of Standards Report 9587 (January 1967)

## BIBLIOGRAPHY

58. Tauber, S.J. et al., "Developing computer programs for searching specific and generic structures, including the formatting of Markush structures", in "Progress in techniques for manipulating and organising chemical information", pp. 23-34 National Bureau of Standards Report 9587 (January 1967)
59. Fraser Williams (Scientific Systems) Ltd., "Proposals for the use of CROSSBOW techniques in the computer handling of Markush structures", Unpublished internal document (1978)
60. Krishnamurthy, E.V., Sankar, P.V., Krishnan, S., "ALWIN - Algorithmic Wiswesser Notation system for organic compounds", Journal of Chemical Documentation **14** (3) 130-141 (1974)
61. Krishnan, S., Krishnamurthy, E.V., "Compact grammar for Algorithmic Wiswesser Notation using Morgan name", Information Processing and Management **12** 19-34 (1976)
62. Krishnan, S., "Compact ALWIN grammar and algorithms for computer handling of organic chemical structures", Ph.D. Thesis, Molecular Biophysics Unit, Indian Institute of Science, Bangalore, India (1975)
63. Howe, W.J., Hagadone, T.R., "Progress toward an online chemical and biological information system at the Upjohn company", in "Retrieval of Medicinal Chemical Information", eds. Howe, W.J., Milne, M.M., Pennel, A.F. ACS Symposium Series **84** 107-131 (1978)
64. Howe, W.J., Hagadone, T.R., "Molecular substructure searching: computer graphics and query entry methodology", Journal of Chemical Information and Computer Sciences in press (1982)
65. Lynch, M.F., "Screening large chemical files", in "Chemical Information Systems", eds. Ash, J.E., Hyde, E., pp. 177-194 Ellis Horwood (1974)
66. Ash, J.E., "Connection Tables and their role in a system", in "Chemical Information Systems", ed. Ash, J.E., Hyde, pp. 156-176, Ellis Horwood (1975)
67. Welford, S.M., "Topological Chemical Grammars and the Generation of Limited Environment Fragments from Generic Chemical Structures", Unpublished Ph.D. Thesis, University of Sheffield (1982)
68. Hopcroft, J.E., Ullman, J.D., "Formal languages and their relation to automata", Addison-Wesley (1969)
69. Ginsburg, S., "The Mathematical Theory of Context-Free Languages", McGraw Hill, New York (1966)

## BIBLIOGRAPHY

70. Hopcroft, J.E., Ullman, J.D., "Introduction to Automata Theory, Languages and Computation" Addison-Wesley, 1979
71. Backhouse, R.C., "Syntax of Programming Languages. Theory and Practice", Prentice Hall (International Series in Computer Science) (1979)
72. Cleaveland, J.C., Uzgalis, R.C., "Grammars for Programming Languages", Elsevier North-Holland (1977)
73. Salomaa, A., "Formal Languages", Academic Press (1973)
74. Chomsky, N., "On certain formal properties of grammars", Information and Control 2 137-167 (1959)
75. Chomsky, N., "Syntactic Structures", Mouton & Co., 's-Gravenhage (1957)
76. Chomsky, N., Miller, G.A., "Introduction to the formal analysis of natural languages", in "Handbook of Mathematical Psychology" eds Luce, R.D., Bush, R.R. and Galanter, E., Volume 2, pp. 269-321, Wiley (1963)
77. Chomsky, N., "Formal properties of grammars", in "Handbook of Mathematical Psychology", eds Luce, R.D. Bush, R.R. and Galanter E., Volume 2, pp. 323-418, Wiley (1963)
78. Chomsky, N., "Aspects of the theory of syntax", MIT Press (1965)
79. Ginsburg, S., Greibach, S.A., "Deterministic context-free languages", Information and Control 9 (6) 620-648 (1966)
80. Aho, A.V., Ullman, J.D., "The Theory of Parsing, Translation and Compiling", 2 Vols, Prentice Hall (1972)
81. Aho, A.V., Ullman, J.D., "Principles of Compiler Design", Addison-Wesley (1977)
82. Knuth, D.E., "On the translation of languages from left to right", Information and Control, 8 607-639 (1965)
83. Lewis, P.M., Stearns, R.E., "Syntax-directed transduction", Journal of the Association for Computing Machinery 15 (3) 465-488 (1968)
84. Rosenkrantz, D.J., Stearns, R.E., "Properties of deterministic top-down grammars", Information and Control 17 226-256 (1970)

BIBLIOGRAPHY

85. Greibach, S., "A new Normal-Form theorem for context-free phrase structure grammars", Journal of the Association for Computing Machinery 12 42-52 (1965)
86. Koranjak, A.J., Hopcroft, J.E., "Simple deterministic languages", Proceedings 7th Annual Symposium on Switching and Automata Theory, IEEE Conference Record, Publication No. 16-C-40 pp. 36-46 (1966)
87. Beatty, J.C., "On the relationship between LL(1) and LR(1) grammars", Journal of the Association for Computing Machinery 29 (4) 1007-1022 (1982)
88. Knuth, D.E., "Top-down syntax analysis", Acta Informatica 1 79-110 (1971)
89. Foster, J.M., "Automatic Syntax Analysis", Macdonald (1970)
90. Sammet, J.E., "Programming Languages. History and Fundamentals", Prentice Hall (1969)
91. Burkhardt, W.H., "Metalanguage and syntax specification", Communications of the Association for Computing Machinery 8 (5) 304-5 (1965)
92. Naur, P., "Report on the algorithmic language ALGOL60", Communications of the Association for Computing Machinery 3 (5) 299-314 (1960)
93. Naur, P. (ed.), "Revised report on the algorithmic language ALGOL60", Communications of the Association for Machinery 6 (1) 1-20 (1963); Computer Journal 5 349-367 (1963); Numerische Mathematik 4 420-452 (1963)
94. Woodward, P.M., Bond, S.G., "User's Guide to Algol68R", Her Majesty's Stationery Office, London (1974)
95. Jensen, K., Wirth, N., "Pascal User Manual and Report", Springer Verlag, New York (1975)
96. Pyle, K., "The Ada Programming Language", Prentice Hall (1981)
97. Knuth, D.E., "Backus Normal Form vs. Backus Naur Form" (letter to the editor), Communications of the Association for Computing Machinery 7 (12) 735-6 (1964)
98. Ginsburg, S., Rice, H.G., "Two families of languages related to ALGOL", Journal of the Association for Computing Machinery 9 350-371 (1962)
99. Iverson, K.E., "A method of syntax specification", Communications of the Association for Computing

## BIBLIOGRAPHY

- Machinery 7 (10) 588-9 (1964)
100. Taylor, W., Turner, L., Waychoff, R., "A syntactical chart of ALGOL 60", Communications of the Association for Computing Machinery 4 393 (1961)
  101. Wirth, N., "What can we do about the unnecessary diversity of notations for syntactic definitions?", Communications of the Association for Computing Machinery 20 (11) 822-3 (1977)
  102. Irons, E.T., "A syntax directed compiler for ALGOL 60", Communications of the Association for Computing Machinery 4 51-55 (1961)
  103. Conway, M.E., "Design of a separable transition-diagram compiler", Communications of the Association for Computing Machinery 6 (7) 396-408 (1963)
  104. Pratt, T.W., "Programming languages: design and implementation", Prentice Hall, Englewood Cliffs, N.J. (1975)
  105. Brown, P.J., "Writing Interactive Compilers and Interpreters", Wiley (1979)
  106. Nicholls, J.E., "The Structure and Design of Programming Languages", Addison-Wesley (1975)
  107. Floyd, R.W., "The syntax of programming languages - a survey", IEEE Transactions on Electronic Computers 13 (4) 346-53 (1964)
  108. Tennent, R.D., "Principles of Programming Languages", Prentice Hall (International Series in Computer Science) (1981)
  109. Wirth, N., "Algorithms + Data Structures = Programs", Prentice Hall (1976)
  110. Welsh, J., McKeag, M., "Structured System Programming", Prentice Hall (International Series in Computer Science) (1980)
  111. Wirth, N., "The programming language Pascal", Acta Informatica 1 35-63 (1971)
  112. Addyman, A.M., "A draft proposal for Pascal", SIGPLAN Notices 15 (4) 1-69 (1980)
  113. Addyman, A.M., et al. "A draft description of Pascal", Software Practice and Experience 9 (5) 381-424 (1979)
  114. Wilson, I.R., Addyman, A.M., "A Practical

## BIBLIOGRAPHY

- Introduction to Pascal", Macmillan, London (1978)
115. Findlay, W., Watt, D.A., "Pascal: an Introduction to Methodical Programming", Pitman (1978)
  116. Atkinson, L.V., "Pascal Programming", Wiley (1980)
  117. McGregor, J.J., Watt, A.H., "Simple Pascal", Pitman Books (1981)
  118. Wirth, N., "The design of a Pascal compiler", Software Practice and Experience 1 309-333 (1971)
  119. Stevens, R., Graham, I., "Matching the abilities of man and computer", Practical Computing pp. 63-7 (August 1979)
  120. Fletcher, D., "Pascal Power", Datamation 25 (8) 142-5 (1979)
  121. Tenenbaum, A.M., Augenstein, M.J., "Data Structures using Pascal", Prentice Hall (1981)
  122. Lecarme, O., Desjardins P., "More comments on the programming language Pascal", Acta Informatica 4 231-243 (1975)
  123. Habermann, A.N., "Critical comments on the programming language Pascal", Acta Informatica 3 47-57 (1973)
  124. Conradi R., "Further critical comments on Pascal, particularly as a systems programming language", SIGPLAN Notices 11 (11) 8-25 (1976)
  125. Main, A., "Pascal creates some bad habits", Computer Weekly (770), p.11 (13 Aug 1981)
  126. Welsh, J., Sneeringer, W.J., Hoare, C.A.R., "Ambiguities and insecurities in Pascal", Software Practice and Experience 7 685-696 (1977)
  127. Cole, J., Morrison, R., "Linguistic Disfigurement", Computing 8 (4) 45 (30 Oct 1980)
  128. Tanenbaum, A.S., "A comparison of Pascal and Algol 68", Computer Journal 21 (4) 316-323 (1978)
  129. Tennent, R.D., "A note on files in Pascal", BIT 17 362-6 (1977)
  130. Doty, K.L., "A top-down evaluation of Pascal", Computer Design 19 (5) 167-177 (1980)
  131. Wirth, N., "An assessment of the programming language

## BIBLIOGRAPHY

- Pascal", IEEE Transactions on Software Engineering SE-1 (2) 192-198 (1975)
132. Glass, R.L., "From Pascal to Pebbleman... and beyond", Datamation 25 (8) 146-50 (1979)
133. Pyle, I., "Ada proves it is equal to 'almost impossible' feats", Computing 9 (44) 26-27 (29 October 1981)
134. Hoare, C.A.R., "The Emperor's Old Clothes" (1981 ACM Turing Award Lecture), Communications of the Association for Computing Machinery 24 75-83 (1981)
135. Skelly, P.G., "The ACM position on standardisation of the Ada language", Communications of the Association for Computing Machinery 25 (2) 118-120 (1982)
136. Ledgard, H.F., Singer, A., "Scaling down Ada (or, Towards a standard Ada subset)", Communications of the Association for Computing Machinery 25 (2) 121-125 (1982)
137. Ince, D., "Why ADA shouldn't go to war", The Guardian p.13 (25th June 1981)
138. Gilbert, J.R., Martin, C.W., "Sheffield University V-mode Pascal System for PR1ME Computers", Unpublished Information File, Computing Services Department, University of Sheffield (1981)
139. de Bakker, J.W., "Semantics of programming languages", Advances in Information Systems Science 2 173-227 (1969)
140. Pagan, F.G., "Formal Specification of Programming Languages", Prentice Hall (1981)
141. Milne, R., Strachey, C., "A Theory of Programming Language Semantics", Chapman and Hall, London (1976)
142. Hoare, C.A.R., Wirth, N., "An axiomatic definition of the programming language Pascal", Acta Informatica 2 335-355 (1973)
143. Hoare, C.A.R., "An axiomatic basis for computer programming", Communications of the Association for Computing Machinery 12 (10) 576-80, 583 (1969)
144. Kupka, I., Wilsing N., "Conversational Languages", Wiley (1980)
145. "The Interpretive Basic Programmer's Guide", Prime Computer Inc., IDR 1813 (1978)

## BIBLIOGRAPHY

146. Spencer, D.D., "A Guide to BASIC Programming; a Time-Sharing Language", Addison-Wesley (1970)
147. Barron, D.W., "Approaches to conversational Fortran", Computer Journal 14 (2) 123-127 (1971)
148. Atkinson, L.V., McGregor, J.J., "CONA - a conversational Algol system", Software - Practice and Experience 8 699-708 (1978)
149. Atkinson, L.V., North, S.D., "COPAS - a conversational Pascal system", Software: Practice and Experience 11 819-829 (1981)
150. Anon., "Abacus makes Pascal as easy as Basic", Computing 8 (17) 6 (24 Apr 1980)
151. Atkinson, L.V., McGregor, J.J., North, S.D., "Context sensitive editing as an approach to incremental compilation", Computer Journal 24 (3) 222-229 (1981)
152. Henry, W.M., Leigh, J.A., Tedd, L.A., Williams, P.W., "Online Searching. An Introduction", Butterworth (1980)
153. Hall, J.L., "Online Information Retrieval Sourcebook", Aslib (1977)
154. Morgan, M.M., Beaman, P.D., Shusman, D.J., Hupp, J.A., Zielstorff, R.D., Barnett, G.O., "Medical query language", Proceedings of the Annual Symposium on Computer Applications in Medical Care, Fifth Symposium, Washington D.C., 1-4 November 1981, 5 322-325 (1981)
155. Orf, H.V., "CHMTRN - LHASA Data Table Language", Unpublished Document, LHASA Project (July 1978)
156. Corey, E.J., Wipke, W.T., Cramer, R.D., Howe, W.J., "Computer-assisted synthetic analysis. Facile man-machine communication of chemical structure by interactive computer graphics", Journal of the American Chemical Society 94 421-430 (1972)
157. Corey, E.J., Cramer, R.D., Howe, W.J., "Computer-assisted synthetic analysis for complex molecules. Methods and procedures for machine generation of synthetic intermediates", Journal of the American Chemical Society 94 440-459 (1972)
158. Dyott, T.M., "Utilization of Stereochemistry and Other Aspects of Computer-Assisted Synthetic Design", Ph.D. Thesis, Princeton University (1973). University Microfilms Order No. 74-9677. Abstracted in Dissertation Abstracts 34 5382-B
159. Wipke, W.T., Braun, H., Smith, G., Choplin, F.,

## BIBLIOGRAPHY

- Seiber, W., "SECS - Simulation and Evaluation of Chemical Synthesis: strategy and planning", in Computer Assisted Organic Synthesis, eds. Wipke, W.T. Howe, W.J., ACS Symposium Series 61 97-127 (1977)
160. Wipke, W.T., Ouchi, G.I., Krishnan, S., "Simulation and evaluation of chemical synthesis - SECS: an application of artificial intelligence techniques", Artificial Intelligence 11 173-193 (1978)
161. Elder, M., Unpublished work (1981)
162. Lin, C.H., Lee, P.L., Lin, J.L., "A SEFLIN compiler - automatic syntactic analysis of separate feature linear notations of chemical compounds", Proceedings of the National Science Council of the Republic of China. Part A 5 (3) 165-172 (1981)
163. Lin, C.-H., "SEFLIN - separate feature linear notation system for chemical compounds", Journal of Chemical Information and Computer Sciences 18 (1) 41-7 (1978)
164. Fehder, P.L., Burnett, M.P., "Syntactic scanning of chemical information", Journal of Chemical Documentation 5 8-13 (1965)
165. Carpenter, N., "Syntax-directed translation of organic chemical files into their two-dimensional representation", Computers in Chemistry 1 25-28 (1976)
166. Barker, P.G., "Syntactic definition and parsing of molecular formulas", Computer Journal 18 355-359 (1978)
167. Kirby, G.H., Morgan, C.H., Rayner, J.D., "Microcomputer Formulae", Education in Chemistry 15-16 (1981)
168. Garfield, E., "An algorithm for translating chemical names to molecular formulas", Journal of Chemical Documentation 2 177-179 (1962)
169. Dyson, G.M., "A cluster of algorithms relating the nomenclature of organic compounds to their structure matrices and ciphers", Information Storage and Retrieval 2 159-199 (1964)
170. Vander Stouw, G.G., Elliott, P.M., Isenberg, A.C., "Automated conversion of chemical substance names to atom-bond connection tables", Journal of Chemical Documentation 14 185-193 (1974)
171. Tauber, S.J., Rankin, K., "Valid structure diagrams and chemical gibberish", Journal of Chemical Documentation 12 (1) 30-34 (1972)

## BIBLIOGRAPHY

172. Rankin, K., Tauber, S.J., "Linguistics as a basis for analysing chemical structure diagrams", Journal of Chemical Documentation **11** (3) 139-141 (1971)
173. Whitlock, H.W., "An organic chemist's view of formal languages", ACS Symposium No 61 pp. 60-80 (American Chemical Society) (1976)
174. Welford, S.M., Lynch, M.F., Barnard, J.M., "Computer storage and retrieval of generic chemical structures in patents. 3. Chemical grammars and their role in the manipulation of chemical structures", Journal of Chemical Information and Computer Sciences **21** (3) 161-168 (1981)
175. Lynch, M.F., Barnard, J.M., Welford, S.M., "Computer storage and retrieval of generic chemical structures in patents. 1. Introduction and general strategy", Journal of Chemical Information and Computer Sciences **21** (3) 148-150 (1981)
176. Barnard, J.M., Lynch, M.F., Welford, S.M., "Computer storage and retrieval of generic chemical structures in patents. 2. GENSAL, a formal language for the description of generic chemical structures", Journal of Chemical Information and Computer Sciences **21** (3) 151-161 (1981)
177. Barnard, J.M., Lynch, M.F., Welford, S.M., "Computer storage and retrieval of generic chemical structures in patents. 4. An extended connection table representation for generic structures", Journal of Chemical Information and Computer Sciences **22** (3) 160-164 (1982)
178. Hill, S.J., "Coding Manuals for Chemical Structure Systems", Unpublished M.Sc. Dissertation, Sheffield University (1981)
179. Feldmann, R.J., Milne, G.W.A., Heller, S.R., Fein, A., Miller, J.A., Koch, B., "An interactive substructure search system", Journal of Chemical Information and Computer Sciences **17** (3) 157-163 (1977)
180. Elder, M., Hull, S.E., Machin, P.A., Mills, O.S., "CSSR. Crystal Structure Search Retrieval. User Manual", 2nd Edition, Daresbury Laboratory, Science and Engineering Research Council, (1981)
181. Milne, G.W.A., Heller, S.R., "NIH/EPA Chemical Information System", Journal of Chemical Information and Computer Sciences **20** (4) 204-211 (1980)
182. Gluck, D.J., "A chemical structure storage and search system developed at Du Pont", Journal of Chemical Documentation **5** 43-51 (1965)

## BIBLIOGRAPHY

183. Roos-Kozel, B.L., Jorgensen, W.L., "Computer-assisted mechanistic evaluation of organic reactions. 2. Perception of rings, aromaticity and tautomers", Journal of Chemical Information and Computer Sciences **21** (2) 101-111 (1981)
184. Corey, E.J., Wipke, W.T., Cramer, R.D., Howe, W.J., "Techniques for perception by a computer of synthetically significant structural features in complex molecules", Journal of the American Chemical Society **94** 431-439 (1972)
185. Corey, E.J., Petersson, G.A., "An algorithm for machine perception of synthetically significant rings in complex cyclic organic structures", Journal of the American Chemical Society **94** 460-465 (1972)
186. Fugmann, R., Nickelsen, H., Nickelsen, I., Winter, J.H., "Representation of concept relations using the TOSAR system of the IDC. Treatise III on Information Retrieval Theory", Journal of the American Society for Information Science **25** (5) 287-307 (1974)
187. Fugmann, R., Nickelsen, H., Nickelsen, I., Winter, J.H., "TOSAR - a topological method for the representation of synthetic and analytical relations of concepts", Angewante Chemie International Edition **9** (8) 589-95 (1970)
188. Dittmar, P.G., Stobaugh, R.E., Watson, C.E., "The Chemical Abstracts registry system. 1. General design", Journal of Chemical Information and Computer Sciences **16** (2) 111-121 (1976)
189. Mockus, J., Isenberg, A.C., Vander Stouw, G.G., "Algorithmic generation of Chemical Abstracts index names. 1. General design", Journal of Chemical Information and Computer Sciences **21** (4) 183-95 (1981)
190. Hoare, C.A.R., "Recursive data structures", International Journal of Computer and Information Studies, **4** (2) 105-132 (1975)
191. Burton, W., "Generalised recursive data structures", Acta Informatica **12** 95-108 (1979)
192. Kinsella, J.E., "An Interactive Test Editor for GENESIS, a Generic Chemical Structure Information System", Unpublished M.A. Dissertation, University of Sheffield (1982)
193. Nishida, F., Takamatsu, S., "Structured-information extraction from patent-claim sentences", Information Processing and Management **18** (1) 1-13 (1982)

## BIBLIOGRAPHY

194. Frome, E., "Effects of information storage and retrieval techniques and computers on problems of patentability", Journal of Chemical Documentation **6** (2) 66-71 (1966)
195. Blick, A.R., "Computer-assisted chemical synthesis packages: is this a new problem in patentability?", Journal of Information Science **1** 227-229 (1979)