

FORMAL DEVELOPMENT OF CORRECT ALGORITHMS: AN
EXAMPLE BASED ON EARLEY'S RECOGNISER

C B Jones
Programming Technology Dept
Product Test Laboratory
IBM Laboratories Ltd
Hursley Park
WINCHESTER, Hants

ABSTRACT

This paper contains the formal development of a correct algorithm from an implicit definition of the task to be performed. Each step of the development can be accompanied by a proof of its correctness. As well as ensuring the correctness of the final program, the structured development gives considerable insight into the algorithm and possible alternatives. The example used is a simplified form of the recognition algorithm due to Earley.

0 INTRODUCTION

The published literature on proving programs correct has tended to confine itself to rather simple algorithms. Apart from the obvious difficulties of long papers, part of the reason for avoiding proofs of larger programs may be that they are almost invariably incorrect! In order to avoid writing a plausible algorithm and then constructing an equally plausible proof which misses the errors, it is necessary for the proof construction to be extremely detailed. The combination of large programs and detailed proofs makes the whole picture difficult to grasp: any errors that are present become elusive.

Is there an alternative? This paper supports an affirmative answer by going through a formal development of a reasonable-sized algorithm. The steps of development represent a, hopefully, natural evolution from the problem specification to the final program. Each step is stated in a formal notation which makes it possible to construct the proof.† Although the overall size may appear large, the structure of the development makes it easy to absorb.

At each stage of development certain assumptions are made which may range from implicit function definitions to properties of data types and their operations. The same stage will, in turn, be further defining functions or adding properties of data types. A proof can be constructed that, under the new set of assumptions, this stage of development fulfils the earlier assumptions. The new assumptions provide, of course, the specification for the next stage. Thus, at each level, it is only necessary to prove that the hypotheses of the preceding level hold: not to prove again earlier results. Moreover, the development is not a chain of equivalences. It is a completely structured development in which it is possible to think, and prove theorems, about one set of problems at a time.

An example of this development taken from the paper begins with an implicit definition of a set (i.e. an object with no order); at the next stage operations are given whose closure on an initial set is shown to satisfy the implicit definition; a mapping of the set onto a list is then given as the next stage and the problems of mimicking the closure are considered.

Although the author ventures some general comments in the next section, this paper does not purport to present a worked out method of formal development. The hope is to encourage further work on a range of examples. Such work should throw light on some of the problems referred to in the discussion of Section 10. In particular there are some areas requiring a much more thorough approach. However, if this is handled correctly further developments should be less difficult than that given below because they can appeal to general theorems.

The body of the paper (Sections 3 to 7 and 9), consists of the development of the algorithm. The example chosen was a recogniser using a simplified form of the ideas of Earley. The reader is asked to refer to ref 1 for a description since such a readable presentation would make detailed exposition, in the current paper, pointless. Suffice it to say, that the algorithm carries out all possible top down parses in parallel, keeping track of them in sets of states. (The simplification made for the current paper was to omit Earley's H_k look-ahead.)

As a result of a draft of this paper the question of modifying developed programs was posed: an attempt to address this question is made in Section 8 with an, albeit rather simple, modification.

† The proofs, which are omitted from the current version of the paper for reasons of space, are to be published in a report available from the author.

1 GENERAL COMMENTS

As already mentioned, this section is far from presenting a complete general method for Formal Development. It is confined to making some observations resulting from the development contained below.

The most pervasive rule in Formal Development appears to be to define as little as possible at each step. Thus only enough new properties are introduced to give some point to the step. Addition of too much, firstly, clouds the structure of the proof and, secondly, may be regretted. The process is addition of properties at each step, removal would require backing up to the point of introduction. An example of the addition of properties is the use of (finite) sets as data types until the main lines of the algorithm are clear, at which time the additional properties of ordering are considered.

The actual mode of development used by the author was at each step to sketch a solution using "understood" notation; then to go back and, in proving that step formally, write down exactly what properties were required of the notation. This provides a detailed specification for the next level.

There is an important division in the subsequent development between the areas with which the development is concerned. Sections 4 and 5 are exclusively related to the Problem domain, the subsequent development is related more to the Data and Language domains.

Comments on the process of modification are left to Section 8.

2 NOTATION

This section reviews the notation used below. Whilst an attempt has been made to minimize the quantity of non-standard notation, use has been made of parts of, so called, VDL (see ref 2 for a fuller treatment). Those parts adopted are defined below in a way which, hopefully, makes this paper self contained.

normal function notation: conditional expressions, <u>where</u> clauses, etc.	standard meanings
(2-1) $f:A \rightarrow B$	f is total over A
$f_a : B \rightarrow C$	specialisation of $f:A \times B \rightarrow C$ to $a \in A$
$\in, \notin, \subseteq, \cup, \cap, \phi, \{x p(x)\}$	normal set operations
$\beta(P) = \{s s \subseteq P\}$	the power set (i.e. set of all subsets)
$N = \{1, 2, 3, \dots\}$	natural numbers
$N^0 = \{0, 1, 2, \dots\}$	integers
$N^n = \{1, 2, \dots, n\}$	natural numbers up to n
$\sim, \vee, \wedge, \supset, \equiv, \exists$	normal logical connectives
$(\exists x) (p(x))$ $x \in X$	bounded quantification
$(\exists! x) (p(x))$	there exists exactly one x ,
$(\iota x) (p(x))$	i.e. $(\exists x) (p(x)) \wedge (p(x) \wedge p(y) \supset x=y)$ that unique x , under above hypothesis
(2-2) $a \& b$	the conditional expression form of "and" (i.e. false if first operand yields false even if second is undefined)

Certain sets of elementary objects are given (e.g. N). Composite objects can be viewed as finite tree structures with elementary objects at the terminal nodes. Branches are named with, so called, selectors. No two branches emanating from the same node may have the same selector name. Subtrees of a given object are called components. As well as naming branches, selectors are used to select components from given objects.

$s(x)$	the component of x named s
(2-3) $is-p = \langle s-1:is-p-1 \rangle, \langle s-2:is-p-2 \rangle, \dots, \langle s-n:is-p-n \rangle$	predicate defining a class of objects each with n sub-components etc.
(2-4) $P = \{x is-p(x)\}$	set of objects satisfying $is-p$
(2-5) $s-1:P \rightarrow P-1$ $s-2:P \rightarrow P-2$ \vdots $s-n:P \rightarrow P-n$	selectors of 2-3 viewed as functions
$s-2(s-1(x))$	iterated selection
Selectors are used in the sequel only when such a component exists. The convention of using the prefixes "is-" and "s-", for predicates and selectors respectively, is followed. Certain "obvious extensions" of strict VDL are employed, e.g.	
(2-6) $is-p = is-\langle s-1:is-p-1 \rangle, \langle s-2:is-p-2 \rangle$	
(2-7) $L = [e1-1, e1-2, \dots, e1-n]$	lists can be considered as objects whose components are named by a subset of the natural numbers (N^n).
(2-8) $\ell(L) = n$	length of list
(2-9) for $1 \leq i \leq \ell(L): L_i = e1-i$ $is-p-list$	selection shown by conventional subscript notation class of all finite lists whose elements satisfy $is-p$
(2-10) $\ell:P-LIST \rightarrow N^0$	
(2-11) $N^n:P-LIST \rightarrow P$	
(2-12) $is-p-list^0$	used for lists which index from zero
(2-13) $is-p-set$	denotes a class of finite sets whose members satisfy $is-p$
(2-14) $\in:P \beta(P) \rightarrow \{T, F\}$	"is a member of" written as an infix predicate symbol
(2-15) $x \in \{e p(e)\} \equiv p(x)$	implicit set definition
(2-16) $U:S \times S \rightarrow S$	set union, such that $x \in (A \cup B) \equiv (x \in A \vee x \in B)$
(2-17) elements of $is-p-set$ cannot be "duplicates"	

3 PROBLEM DEFINITION

This section gives the definition of the task to be performed by the algorithm developed in subsequent sections. In outline: the class of grammars is defined in 3-1 to 3-4 in the familiar terminal/non-terminal way; the class of strings which can be produced from a given grammar is defined in 3-5 and 3-6; a function "root" is introduced in 3-9 which is assumed to provide the non-terminal from which acceptable strings must be derivable. The recognition task, for given string and grammar, can now be defined in 3-10 using the above notions. A useful Lemma on productions is given in 3-7.

The class of grammars is defined:

$$(3-1) \quad \text{is-grammar} = \text{is-rule-set}$$

which is assumed to have properties 2-13, 2-14

$$(3-2) \quad \text{is-rule} = (\langle \text{s-lhs} : \text{is-nt} \rangle, \langle \text{s-rhs} : \text{is-el-list} \rangle)$$

properties 2-5 and 2-10, 2-11 are assumed

\bar{r} will be used as an abbreviation for $\ell(\text{s-rhs}(r))$

$$(3-3) \quad \text{is-el} = \text{is-nt} \vee \text{is-t}$$

$$(3-4) \quad \text{NT} \cap \text{T} = \emptyset \quad \text{non-terminal and terminal sets are disjoint}$$

The following production symbols are defined:

for $\text{is-el-list}(\alpha)$, $\text{is-el-list}(\beta)$, $\text{is-grammar}(G)$:

$$(3-5) \quad \alpha A \beta \stackrel{*}{\Rightarrow}_G \alpha \gamma \beta \quad \equiv (\exists r_{r \in G}) (\text{s-lhs}(r) = A \wedge \text{s-rhs}(r) = \gamma)$$

$$(3-6) \quad \alpha \stackrel{*}{\Rightarrow}_G \beta \quad \equiv (\exists m) (\exists \alpha_0, \alpha_1, \dots, \alpha_m) (\alpha = \alpha_0 \Rightarrow \alpha_1 \Rightarrow \alpha_2 \Rightarrow \dots \Rightarrow \alpha_m = \beta)$$

In the sequel, the grammar can be omitted from the production symbol with no danger of confusion.

The following Lemma can be proven from 3-5, 3-6.

$$(3-7) \quad \text{Lemma} \quad \gamma \stackrel{*}{\Rightarrow} \tau \supset (\exists t_1, \dots, t_{\ell(\gamma)}) (1 \leq t_1 \leq t_2 \leq \dots \leq t_{\ell(\gamma)} = \ell(\tau) \wedge \\ \gamma_1 \stackrel{*}{\Rightarrow} \tau_1 \dots \tau_{t_1} \wedge \\ \vdots \\ \gamma_{\ell(\gamma)} \stackrel{*}{\Rightarrow} \tau_{t_{\ell(\gamma)}-1+1} \dots \tau_{t_{\ell(\gamma)}})$$

The "specification" of the algorithm to be developed, say rec , can now be stated as follows:

$$(3-8) \quad \text{rec} : \text{T-LIST} \times \text{GRAMMAR} \rightarrow \{\text{YES}, \text{NO}\}$$

$$(3-9) \quad \text{such that for is-t-list} \quad \text{assuming properties 2-10, 2-11}$$

$$\text{is-grammar}(G) \quad \text{see 3-1}$$

$$\text{root} : \text{GRAMMAR} \rightarrow \text{NT} \quad \text{a predefined function}$$

$$(3-10) \quad \text{rec}(X, G) = \text{YES} \quad \text{iff } \text{root}(G) \stackrel{*}{\Rightarrow} X \\ \text{NO} \quad \text{otherwise}$$

Given that 3-8 requires that rec be total (see 2-1), the second part of 3-10 is not proved explicitly.

4 EARLEY'S STATE SETS

The reader is assumed to be familiar with the concepts of states and state sets as described by Earley in ref 1. (The simplification in the current paper drops the fourth element of Earley's states). Collections of state sets correspond to objects satisfying is-S (see 4-2) from which individual state sets can be retrieved using the function "state-set". The only properties required of the data object containing state-sets (see 4-3) are the applicability of the predicate "is member of" and the lack of duplicate members. Objects containing states are defined by 4-4 to be triples whose contents will be referred to via given selectors or using the specified abbreviation. More properties of the storage will be added as the algorithm evolves, those given so far are enough to show the first ideas about the algorithm.

It is convenient to assume that there is only one rule of the grammar of which root (G) is the left hand side, this restriction is made in 4-1 since it has no real effect on the class of valid grammars:

$$(4-1) \quad (\exists! r) (s\text{-lhs}(r) = \text{root}(G)) \\ r \in G$$

The object containing all state-sets satisfies:

$$(4-2) \quad \text{is-S}$$

which is characterised by the existence of two functions:

$$\begin{aligned} \text{len} : S &\rightarrow N \\ \text{state-set} : N^n \times S &\rightarrow \text{STATE-SET} \\ &\text{where } N^n \text{ is the set } \{0,1,\dots,\text{len}(S)\} \end{aligned}$$

$$(4-3) \quad \text{is-stateset} = \text{is-state-set}$$

of which 2-14 and 2-17 are required to hold.

$$(4-4) \quad \text{is-state} = (\text{<s-rule:is-rule>}, \text{<s-j:is-n>}, \text{<s-f:is-n>})$$

of which 2-5 is required to hold.

The abbreviation:

$$s = \langle r, j, f \rangle \quad \text{is used for the object for which} \quad \begin{aligned} s\text{-rule}(s) &= r \\ s\text{-j}(s) &= j \\ s\text{-f}(s) &= f \end{aligned}$$

States will be used to record partial parses. The presence of a particular state, say $\langle r, j, f \rangle$, in the i th state set will mean that a string can be derived from the grammar, G , which consists of the first f characters of X (the string to be recognised) followed by left hand side of r , and that the first j elements of the right hand side of r can produce, from G , the string $X_{f+1} \dots X_i$. This property is stated formally in 4-7 below and represents an upper bound for valid state-sets: any state present must possess this property. In order to permit the "rec" algorithm to decide recognition, not all of these states are essential so 4-6 and 4-8 set a lower bound on valid state-sets: the former specifies a start element which must be present; the latter specifies that if $\langle r, j, f \rangle$ is in the i th state set and the next (i.e. $j+1$ th) element of the right hand side of r can produce $X_{i+1} \dots X_m$ then $\langle r, j+1, f \rangle$ must be a member of the m th state-set.

It is fundamental to the style of development proposed by this paper that the freedom, to choose which set between the lower and upper bound will be generated, is left for the time being. The properties given permit the development of the first steps of the algorithm (see 4-5) and the proof that, under the assumptions on the creation of S , it fulfils the task described in Section 3. This freedom is used to show the correctness of an optimisation modification in Section 8 and could be used to validate many different algorithms including the use of Earley style look-ahead. Furthermore, the proof of the actual construction of state sets constructed in Section 5 has been made simpler by only having to prove that these looser properties hold.

The definition of a function, alleged to satisfy 3-8, 3-10, can now be given:

$$(4-5) \quad \text{rec}(X, G) = \begin{array}{l} \langle \overline{rr}, \overline{rr}, 0 \rangle \in \text{state-set}(\ell(x), S) \quad \rightarrow \underline{\text{YES}} \\ T \quad \quad \quad \rightarrow \underline{\text{NO}} \end{array}$$

$\underline{\text{where}} \quad \text{is-} S(S)$
 $\underline{\text{and}} \quad \text{len}(S) = \ell(x)$
 $\underline{\text{and}} \quad \overline{rr} = (\overline{1r})(s\text{-lhs}(r) = \text{root}(G))$
 $\quad \quad \quad r \in G$

The correctness of this function relies on the following assumptions about S . (These assumptions will, of course, form the definition for further development.) Some base element must exist in S :

$$(4-6) \quad \langle \overline{rr}, 0, 0 \rangle \in \text{state-set}(0, S)$$

Any element of a state set must, at least, satisfy the properties:

$$(4-7) \quad 0 \leq i \leq \ell(X) \ \& \ \langle r, j, f \rangle \in \text{state-set}(i, S) \Rightarrow \begin{array}{l} r \in G \ \wedge \\ (\exists \alpha \quad \quad \quad)(\text{root}(G) \stackrel{*}{=} X_1 \dots X_f \text{ s-lhs}(r) \alpha) \wedge \\ \text{is-el-list}(\alpha) \\ (j > 0 \Rightarrow \text{s-rhs}(r)_1 \dots \text{s-rhs}(r)_j \stackrel{*}{=} X_{f+1} \dots X_i) \wedge \\ j \leq \overline{r} \ \wedge \ 0 \leq f \leq i \end{array}$$

Certain elements must be in the state-sets:

$$(4-8) \quad 0 \leq i < \ell(x) \ \& \ \langle r, j, f \rangle \in \text{state-set}(i, S) \wedge j \neq \overline{r} \ \& \ (\exists m \quad \quad \quad)(\text{s-rhs}(r)_{j+1} \stackrel{*}{=} X_{i+1} \dots X_m)$$

$i+1 \leq m \leq \ell(X)$
 $\Rightarrow \langle r, j+1, f \rangle \in \text{state-set}(m, S)$

$$(4-9) \quad S \text{ must be created in a finite amount of time.}$$

Property 4-9 above is given to ensure that any subsequent algorithm generates state-sets in a useful way and permits the proof (see 4-13) that 4-5 is total. After a trivial Lemma (4-10) on the finiteness of state-sets, the proofs that both implications of the first part of 3-10 hold are given in 4-11 and 4-12. Of course, these proofs all rely on the assumptions that S satisfies 4-6 to 4-9.

$$(4-10) \quad \text{Lemma} \quad \text{for } 0 \leq i \leq \text{len}(S) : \text{state-set}(i, S) \text{ is finite.}$$

proof follows from finiteness of G and \overline{r} .

$$(4-11) \quad \text{Theorem} \quad \text{rec}(X, G) = \underline{\text{YES}} \Rightarrow \text{root}(G) \stackrel{*}{=} X$$

proof follows from 4-7.

$$(4-12) \quad \text{Theorem} \quad \text{root}(G) \stackrel{*}{=} X \Rightarrow \text{rec}(X, G) = \underline{\text{YES}}$$

proof follows from 4-6 and 4-8 using 3-7.

$$(4-13) \quad \text{Theorem} \quad \text{rec is total}$$

proof follows from 4-10 and 4-9 and check of domains/ranges.

5 CREATING STATE SETS

Section 4 introduced the concept of state sets and some essential properties thereof, but offered no way of creating them. In fact (s-rhs(rr)) entries, presumably created by magic, would have been sufficient to fulfil the properties stated. This section shows how state-sets can be constructed using the main ideas (i.e. prediction, completion, scanning) of Earley's method. In order to facilitate construction of the required elements, many other states are generated. It is interesting to consider the parallel with proof construction where an induction hypothesis stronger than the theorem statement is used in order to prove the latter. First, the definition of rr which was in the last section a local convention, is adopted for the sequel.

$$(5-0) \quad rr = (\text{lr})(s\text{-lhs}(r) = \text{root}(G)) \quad \text{see 4-1} \\ r \in G$$

Earley's operations of prediction etc, are, of course, described in terms of lists of states. As will be clear shortly, many of the interesting points about these operations can be brought out using sets of states. Again it is found that deferring part of the detail, in this case the special ordering problems imposed by using lists, will structure and clarify the proofs. It will be necessary to use implicit set notation for creation of state sets so now the restriction:

$$(5-1) \quad \text{Property 2-15 holds for is-stateset} \quad \text{is added}$$

Using sets of states which are cumulative, prompts the idea of using the notion of the closure of a set under an operation as generating function. Intuitively the closure of an operation on a set is the minimum set containing the base set and having the property that applying the operation to any element creates only elements already in the set. We shall use the following properties:

$$\text{for} \quad \text{op} : P \rightarrow \beta(P) \quad \text{i.e. from elements of } P \text{ to sets of such elements}$$

$$(5-2) \quad \text{closure} : (op) \times \beta(P) \rightarrow \beta(P) \quad \text{i.e. from sets of elements to sets of elements}$$

$$\text{such that for } S \subseteq P:$$

$$x \in (\text{closure of op on } S) \equiv x \in S \vee (\exists y)_{y \in S} (x \in \text{closure of op on } (op(y)))$$

Notice that closure is total for suitable op. It should be clear that closure is monotonic. That is:

$$(5-3) \quad S_1 \subseteq S_2 \supset (\text{closure of op on } S_1) \subseteq (\text{closure of op on } S_2)$$

Using closure it is now possible to show how the basic operations (i.e. prediction etc) are used to create the statesets. The base element is inserted by 5-4 which also shows that prediction is the only applicable operation on the 0th stateset. The creation, for all other statesets, of a kernel set of states resulting from the scanning operation on the preceding stateset is specified in 5-5. Both the prediction and completion operations are employed on all statesets but the last, where the former is not required. The closure is created from the kernel set described. The importance of the reliance by the completer on S (strictly on a part of S) is returned to in the next section.

$$(5-4)^+ \quad \text{state-set } (0, S) = \text{closure of predict}_{0, X, G} \text{ on } \{<r, 0, 0>\}$$

$$(5-5) \quad \text{for } 1 \leq i \leq l(X): \quad \text{state-set } (i, S)^S = \text{scan } (i, X, G, \text{state-set}(i-1, S))$$

$$(5-6) \quad \text{for } 1 \leq i \leq l(X): \quad \text{state-set } (i, S) = \text{closure of } \left\{ \begin{array}{l} \text{predict}_{i, X, G} \\ \text{complete}_{S, i, X, G} \end{array} \right\} \text{ on state-set } (i, S)^S$$

$$(5-7) \quad \text{state-set } (l(X), S) = \text{closure of complete}_{S, i, X, G} \text{ on state-set } (l(X), S)^S$$

⁺ Notice use of function specialisation, see Section 2.

The definitions of the basic operations of prediction, scanning and completion can now be given. For reasons discussed in Section 9, it is important to guard against predicates becoming undefined by careless use of " \wedge " which is not normally defined for three valued logic. (see ref 4)

$$(5-8) \quad \text{predict}(i, X, G, \langle r, j, f \rangle) = \{ \langle s, 0, i \rangle \mid j \neq \bar{r} \ \& \ \text{is-nt}(s\text{-rhs}(r)_{j+1}) \ \& \ s \in G \ \& \ s\text{-lhs}(s) = s\text{-rhs}(r)_{j+1} \}$$

$$(5-9) \quad \text{scan}(i, X, G, ss) = \{ \langle r, j+1, f \rangle \mid \langle r, j, f \rangle \in ss \wedge j \neq \bar{r} \ \& \ \text{is-t}(s\text{-rhs}(r)_{j+1}) \ \& \ X_i = s\text{-rhs}(r)_{j+1} \}$$

$$(5-10) \quad \text{complete}(S, i, X, G, \langle r, j, f \rangle) = \{ \langle s, m+1, g \rangle \mid j = \bar{r} \ \& \ \langle s, m, g \rangle \in \text{state-set}(f, S) \ \& \ m \neq \bar{s} \ \& \ \text{is-nt}(s\text{-rhs}(s)_{m+1}) \ \& \ s\text{-rhs}(s)_{m+1} = s\text{-lhs}(r) \}$$

In order to resolve any possible difficulties deriving from too heavy a reliance on formalism the example used in ref 1 is now presented in the terms of the current paper. (Notice the addition of R to satisfy 4-1)

$$G = \{r1, r2, r3, r4, r5, r6\} \quad \text{where} \quad \begin{aligned} r1 &= \langle \langle s\text{-lhs} : R \rangle, \langle s\text{-rhs} : [E] \rangle \rangle \\ r2 &= \langle \langle s\text{-lhs} : E \rangle, \langle s\text{-rhs} : [T] \rangle \rangle \\ r3 &= \langle \langle s\text{-lhs} : E \rangle, \langle s\text{-rhs} : [E, +, T] \rangle \rangle \\ r4 &= \langle \langle s\text{-lhs} : T \rangle, \langle s\text{-rhs} : [P] \rangle \rangle \\ r5 &= \langle \langle s\text{-lhs} : T \rangle, \langle s\text{-rhs} : [T, *, P] \rangle \rangle \\ r6 &= \langle \langle s\text{-lhs} : P \rangle, \langle s\text{-rhs} : [a] \rangle \rangle \end{aligned}$$

such that $\text{is-nt}(R), \text{is-nt}(E), \text{is-nt}(T), \text{is-nt}(P), \text{is-t}(\dagger), \text{is-t}(*), \text{is-t}(a), rr = r1$

Consider the recognition of:

$$X = [a, +, a]$$

$$\begin{aligned} \text{state-set}(0, S) &= \text{closure of predict on } \{ \langle r1, 0, 0 \rangle \} \\ &= \text{closure of predict on } \{ \langle r2, 0, 0 \rangle \} \\ &= \text{closure of predict on } \{ \langle r4, 0, 0 \rangle, \langle r5, 0, 0 \rangle \} \\ &= \text{closure of predict on } \{ \langle r6, 0, 0 \rangle \} \\ &= \{ \langle r1, 0, 0 \rangle, \langle r2, 0, 0 \rangle, \langle r4, 0, 0 \rangle, \langle r5, 0, 0 \rangle, \langle r6, 0, 0 \rangle \} \end{aligned}$$

$$\begin{aligned} \text{state-set}(1, S)^S &= \text{scan}(\text{state-set}(0, S)) \\ &= \{ \langle r6, 1, 0 \rangle \} \end{aligned}$$

$$\begin{aligned} \text{state-set}(1, S) &= \text{closure of predict/complete on state-set}(1, S)^S \\ &= \text{closure of predict/complete on } \{ \langle r6, 1, 0 \rangle, \langle r4, 1, 0 \rangle, \langle r5, 1, 0 \rangle \} \\ &= \text{closure of predict/complete on } \{ \langle r2, 1, 0 \rangle, \langle r3, 1, 0 \rangle \} \\ &= \text{closure of predict/complete on } \{ \langle r1, 1, 0 \rangle \} \\ &= \{ \langle r6, 1, 0 \rangle, \langle r4, 1, 0 \rangle, \langle r5, 1, 0 \rangle, \langle r2, 1, 0 \rangle, \langle r3, 1, 0 \rangle, \langle r1, 1, 0 \rangle \} \end{aligned}$$

$$\begin{aligned} \text{state-set}(2, S)^S &= \text{scan}(\text{state-set}(1, S)) \\ &= \{ \langle r3, 2, 0 \rangle \} \end{aligned}$$

$$\begin{aligned} \text{state-set}(2, S) &= \text{closure of predict/complete on state-set}(2, S)^S \\ &= \text{closure of predict/complete on } \{ \langle r3, 2, 0 \rangle, \langle r4, 0, 2 \rangle, \langle r5, 0, 2 \rangle \} \\ &= \text{closure of predict/complete on } \{ \langle r6, 0, 2 \rangle \} \\ &= \{ \langle r3, 2, 0 \rangle, \langle r4, 0, 2 \rangle, \langle r5, 0, 2 \rangle, \langle r6, 0, 2 \rangle \} \end{aligned}$$

$$\begin{aligned} \text{state-set}(3, S) &= \text{scan}(\text{state-set}(2, S)) \\ &= \{ \langle r6, 1, 2 \rangle \} \end{aligned}$$

$$\begin{aligned} \text{state-set}(3, S)^S &= \text{closure of complete on state-set}(3, S)^S \\ &= \text{closure of complete on } \{ \langle r6, 1, 2 \rangle, \langle r4, 1, 2 \rangle \} \\ &= \text{closure of complete on } \{ \langle r3, 3, 0 \rangle, \langle r5, 1, 2 \rangle \} \\ &= \text{closure of complete on } \{ \langle r1, 1, 0 \rangle \} \\ &= \{ \langle r6, 1, 2 \rangle, \langle r4, 1, 2 \rangle, \langle r3, 3, 0 \rangle, \langle r5, 1, 2 \rangle, \langle r1, 1, 0 \rangle \} \end{aligned}$$

Referring to 4-5 gives:

$$\text{rec}(X, G) = \text{YES}$$

The algorithm represents a way of generating some set of states between those specified by 4-8 and 4-7. The justification[†] consists of proving this containment (see 5-12, 5-13) and establishing 4-6 and 4-9 (see 5-11, 5-14 respectively). Notice that this is all! It is not necessary to prove again any of the results dealt with in Section 3.

(5-11) Theorem $\langle rr, 0, 0 \rangle \in \text{state-set}(0, S)$

proof follows from 5-4, 5-3.

(5.12) Theorem

$$0 \leq i < \ell(X) \ \& \ \langle r, j, f \rangle \in \text{state-set}(i, S) \supset r \in G \ \wedge \\ (\exists \alpha) (\text{root}(G) \xRightarrow{*} X_1 \dots X_f \text{ s-lhs}(r) \alpha) \ \wedge \\ \text{is-el-list}(\alpha) \\ (j > 0 \supset \text{s-rhs}(r)_1 \dots \text{s-rhs}(r)_j \xRightarrow{*} X_{f+1} \dots X_i) \ \wedge \\ j \leq \bar{r} \ \wedge \ 0 \leq f \leq i$$

proof: shows that the set satisfying the above is a fixed point of 5-4 to 5-7 (i.e. that given such a set the operations 5-8, 5-9, 5-10 create only elements of that set).

(5-13) Theorem

$$0 \leq i < \ell(X) \ \& \ \langle r, j, f \rangle \in \text{state-set}(i, S) \ \wedge \ j \neq \bar{r} \ \& \ (\exists m)_{i+1 \leq m \leq \ell(X)} (\text{s-rhs}(r)_{j+1} \xRightarrow{*} X_{i+1} \dots X_m) \\ \supset \ \langle r, j+1, f \rangle \in \text{state-set}(m, S)$$

proof by induction on n (the number of \Rightarrow steps in \Rightarrow^* see 3-5, 3-6, using 5-8, 5-9, 5-10.

(5-14) Theorem S is created in finite time

proof follows from the closure and finiteness (4-10)

6 MAPPING STATE SETS ONTO LISTS

So far state sets have been assumed to be set-like objects and they have been extended by closure. Whilst it would be possible to model this state of affairs in a conventional store it would be much more convenient if we could show how the sets can be mapped onto lists. The lists will obviously have extra properties, in particular an ordered property. If closure can now be replaced by a new, more efficient, operation which relies on this ordering, an increase in efficiency results. This section makes precisely this mapping. Subject to a restriction on grammars (see 6-9) the extension is performed by a single scan over the lists. It is, in the author's opinion, worthwhile understanding the cause of the restriction. It would for instance be a relatively simple task, at this level, to avoid introducing the restriction by utilizing a more sophisticated scan.

So far the only properties assumed of state-sets are those required by 4-3 and 5-1. That is:

$$\epsilon: \text{STATE} \times \text{STATASET} \rightarrow \{T, F\} \\ x \in \{\text{state} \mid p(\text{state})\} \equiv p(x) \\ \text{statesets do not contain duplicates}$$

This section describes how these objects can be mapped onto lists. Thus:

(6-1) $\text{is-stateset} = \text{is-state-list}$ assuming properties 2-10 and 2-11.

[†] It should be possible to construct a direct proof that a set given by 4-7, with equivalence replacing implication, is the minimum fixed point (see ref 3) of 5-4 to 5-10. The current author's attempts have foundered on the ordering problems caused by the completer.

Now the membership and implicit set definitions are reinterpreted as follows:

(6-2) $e \in s$ becomes $(\exists_{1 \leq i \leq \ell(s)})(s_i = e)$ but will be abbreviated $e \in s$

(6-3) $\{s \mid p(s)\}$ becomes $\text{create-state-list}(p)$ but will be abbreviated $\{s \mid p(s)\}$ with the following properties:

$$\begin{aligned} \text{let } \ell(\{s \mid p(s)\}) &= n & p(x) &\supset (\exists_{1 \leq j \leq n} !j)(\{s \mid p(s)\}_j = x) \\ & & 1 \leq j \leq n &\supset p(\{s \mid p(s)\}_j) \end{aligned}$$

A new operation will be required to manipulate state-sets.

(6-4) combine $(\text{list}_1, \text{list}_2)$ which will be written \cup

with the following properties:

$$\begin{aligned} s \in (\text{list}_1 \cup \text{list}_2) &\equiv s \in \text{list}_1 \vee s \in \text{list}_2 \\ (\exists_{1 \leq j, k \leq \ell(\text{list}_1 \cup \text{list}_2)} k) ((\text{list}_1 \cup \text{list}_2)_j = (\text{list}_1 \cup \text{list}_2)_k \supset j = k) \\ 1 \leq i \leq \ell(\text{list}_1) &\supset (\text{list}_1 \cup \text{list}_2)_i = (\text{list}_1)_i \end{aligned}$$

The properties required of state-sets can now be verified:

(6-5) Lemma

- a) $\in : \text{STATE} \times \text{STATASET} \rightarrow \{T, F\}$
- b) $x \in \{\text{state} \mid p(\text{state})\} \equiv p(x)$
- c) statesets contain no duplicates

Operating on lists makes it desirable to replace the closure operations (see 5-2) with linear scans. It will be found necessary to introduce a restriction on grammars in order to fulfil property 5-2.

The construction of S is now given by:

(6-6) $S = S^Z$ where next $(z+1) = \text{end}$

(6-7) next $(1) = (0, 1)$
for $i > 1$:

$$\text{next } (i) = \ell(\text{state-set}(j, S^{i-1})) = k \rightarrow j \neq \ell(X) \ \& \ \text{state-set}(j+1, S^{i-1}) \neq [] \rightarrow (j+1, 1)$$

T -> end

T -> (j, k+1)

where $(j, k) = \text{next } (i-1)$

The preceding section represented a solution to the main parts of the recognition task. However, for a number of reasons it is not a program which can be run on a computer. The principal point is that the algorithm of Section 6 uses a storage object which is not representable in any known language, for example the grammar is just "referred to" whenever required. It might be possible to sit at a console and supply rules whenever needed, but a more useful program will result if the first step is to read in both grammar and string to be recognised in such a way that subsequent references are simple. The particular storage structure given below has been adopted from an actual program written in the conventional way and designed with a view to efficiency. This decision was taken so as to avoid the danger of choosing a data structure particularly suited to the proof but which put limitations on the potential efficiency of the final algorithm.

The input of grammar and string will be assumed to be carried out by routines called INGR and INSTR respectively. The grammar is stored in the s-RULES, s-NONT, s-RULED, s-STR and s-STRCHAR; the string in the s-INPUT components of the store given in 7-1 to 7-7 below. What were, in the algorithm references to the original arguments, now become references to these storage components and the reinterpretations to be made are specified in 7-9 to 7-17.

The object is-S is also refined in this section and now occupies the s-S and s-STATE parts of 7-1 to 7-7.

The store is still presented as a VDL style object and a further transition to PL/I Data Structures will be made later. The decision to introduce this step is perhaps questionable, but is probably justified since the object of 7-1 could be realised in other languages. The only properties required are those in 7-8 and although "bunches" were used with PL/I based storage in mind, the only objection to an array style implementation with numeric "selectors" would be that all elements would have to be of the same length.

Store is considered as an object satisfying:

```
(7-1)  is-ξ = (<s-S : is-state-ptr-list0>,
            <s-STATE : is-STATE-bunch>,
            <s-INPUT : is-t-list>,
            <s-RULES : is-(<s-TYPE : is-T ∨ is-N>,
                          <s-N : is-n>) - list>,
            <s-NONT : is-ruled-ptr-list>,
            <s-RULED : is-RULED-bunch>,
            <s-STR : is-(<s-START : is-n>)-list>,
            <s-STRCHAR : is-char-list>)

    where -
(7-2)  is-state-ptr = is-sel
(7-3)  is-STATE-bunch = {<κ : state> || is-sel(κ) ∧ is-STATE(state)}
(7-4)  is-STATE = (<s-SIZE : is-n>,
                  <s-N : is-n>,
                  <s-INFO : is-(<s-RULE : is-(<s-RPTR : is-ruled-ptr>,
                                         <s-SSC : is-n>)>,
                               <s-RULPOS : is-n>,
                               <s-STRPOS : is-n>)-list>)
```

```
(7-5)  is-ruled-ptr = is-sel
(7-6)  is-RULED-bunch = {<κ : ruled> || is-sel(κ) ∧ is-RULED(ruled)}
(7-7)  is-RULED = (<s-NEXT : is-n>,
                  <s-RULED : is-(<s-START : is-n>,
                               <s-LENGT : is-n>)-list>)
```

```
(7-8)  In all object of 7-1 to 7-7:      selectors are assumed to have property 2-5
                                           lists are assumed to have properties 2-10, 2-11
```

Using the above data object as storage the algorithm can now use the following ways of referring to its arguments:

```
(7-9)  ℓ(X)      becomes ℓ(s-INPUT(ξ))
```

```
(7-10) Xi      becomes s-INPUT(ξ)i
```

Now, INGR defines a relation between rules of is-grammar and objects satisfying $is-(\langle s-RPTR : is-ruled-ptr \rangle, \langle s-SSC : is-n \rangle)$ which is written:

$$\begin{aligned}
 & I(r, (p, i)) \\
 (7-11) \quad & \left. \begin{array}{l} \bar{r} \\ \ell(s-rhs(r)) \end{array} \right\} \text{ becomes } s-LENGT(s-RULE(p(s-RULED(\xi)))_i) \\
 (7-12) \quad & s-rhs(r)_j \text{ becomes } s-RULES(\xi) \quad s-START(s-RULED(p(s-RULED(\xi)))_i) + j - 1
 \end{aligned}$$

for convenience, let $I(s-rhs(r)_j)$ be an abbreviation for $s-RULES(\xi) \quad s-START(s-RULED(p(s-RULED(\xi)))_i) + j - 1$

$$(7-13) \quad is-t(s-rhs(r)_j) \text{ becomes } s-TYPE(I(s-rhs(r)_j)) = T$$

$$(7-14) \quad is-nt(s-rhs(r)_j) \text{ becomes } s-TYPE(I(s-rhs(r)_j)) = N$$

$$(7-15) \quad \text{for } is-t(s-rhs(r)_j):$$

$$s-rhs(r)_j \text{ becomes } s-STRCHAR(\xi)_{s-N(I(s-rhs(r)_j))}$$

$$(7-16) \quad \text{for } is-nt(s-rhs(r)_j):$$

$$s-rhs(r)_j = s-lhs(s) \text{ becomes } s-NONT(\xi)_{s-N(I(s-rhs(r)_j))} = q \quad \text{where } I(s, (q, k))$$

$$(7-17) \quad r \in G \text{ becomes } is-RULED(p(s-RULED(\xi))) \wedge i < s-NEXT(p(s-RULED(\xi)))$$

Now, it is necessary to show that the reinterpretations of these functions still satisfy the stated properties (3-8 and 3-1 to 3-4).

(7-18) Lemma

$$\begin{aligned}
 \ell: \quad & T-LIST \rightarrow N^0 && \text{see 7-9, 3-8, 2-10} \\
 N^n: \quad & T-LIST \rightarrow T && \text{see 7-10, 3-8, 2-11}
 \end{aligned}$$

Both results follow immediately from 7-1, 7-8, 2-10, 2-11.

(7-19) Lemma

$$\begin{aligned}
 a) \quad & s-RULES, s-NONT, s-RULED, s-STR, s-STRCHAR \text{ are finite} && \text{see 3-1, 2-13} \\
 b) \quad & \epsilon: (RULED-PTR \times N) \times (RULES, \dots, STRCHAR) \rightarrow \{T, F\} && \text{see 7-17, 3-1, 2-14} \\
 c) \quad & \ell: (RULED-PTR \times N) \times (RULES, \dots, STRCHAR) \rightarrow N && \text{see 7-11, 3-2, 2-10} \\
 d) \quad & N^n: (RULED-PTR \times N) \times (RULES, \dots, STRCHAR) \rightarrow && \text{see 7-15, 7-16, 3-2, 2-11} \\
 e) \quad & \sim(is-nt(s-rhs(r)_j) \wedge is-t(s-rhs(r)_j)) && \text{see 7-13, 7-14, 3-3, 3-4}
 \end{aligned}$$

proofs follow from definitions and 7-8.

However, it is not possible to simply change the problem definition (see 3-10 and 3-5, 3-6,) which is stated in terms of the original operations (e.g. $r \in G$). It was clear from the beginning that routines like INSTR and INGR would be required and it would have been possible to use different operations in the algorithm and note the equivalences which must hold between the original and algorithm notations. This was not done since it would have made the subsequent writing more tedious. In making the change to our re-interpreted operations we must now, carefully, check what properties are required and use them as a specification of the INSTR and INGR routines.

(7-20) Constraints on INSTR -

$$\begin{aligned}
 \ell(X) &= \ell(s-INPUT(\xi)) \\
 X_i &= s-INPUT(\xi)_i
 \end{aligned}$$

(7-21) Constraints on INGR -

```

let (I(r,(p,i))
    r ∈ G           ≡ is-RULED(p(s-RULED(ξ))) ∧ i < s-NEXT(p(s-RULED(ξ)))
     $\bar{r}$               = s-LENGT (s-RULE(p(s-RULED(ξ)))i)
    is-t(s-rhs(r)j) ≡ s-TYPE (I(s-rhs(r)j)) = T
    is-nt(s-rhs(r)j) ≡ s-TYPE (I(s-rhs(r)j)) = N

for is-t(s-rhs(r)j):
    s-rhs(r)j      ≡ s-STRCHAR(ξ)s-N (I(s-rhs(r)j))

for is-nt(s-rhs(r)j):
    s-rhs(r)j = s-lhs(s) ≡ s-NONT(ξ)s-N (I(s-rhs(r)j)) = q
                                where I(s,(q,k))

```

The storage structure given in 7-1 to 7-4 also gives us an interpretation of the collection of state-sets (see 4-2).

```

(7-22) len(S)          becomes     $\ell^0$  (s-S(ξ))
(7-23) state-set(i,S)  becomes    s-S(ξ)i (s-STATE(ξ))
(7-24) for is-state(s): s-rule(s) becomes    s-RULE(s)
(7-25)                  s-j(s)    becomes    s-RULEPOS(s)
(7-26)                  s-f(s)    becomes    s-STRPOS(s)

```

It is, of course, necessary to check that the required properties still hold:

(7-27) Lemma ℓ^0 , state-set, s-RULE, s-RULPOS, s-STRPOS behave as in 4-2, 4-4
 proofs all follow from 7-1, 7-2, 7-3, 7-4, 7-8.

The resulting reinterpretations of is-stateset are:

```

(7-28)  $\ell$ : STATE-SET → N          becomes    s-N(STATE)
(7-29)  $N^n$ : STATE-SET → STATE    becomes    s-INFO(STATE)i
(7-30) Lemma

```

The obvious equivalences hold.

Follows from 7-1, 7-4, 7-8

Notice that $\underline{\leq}$ and $\underline{\cup}$ (see 6-2 and 6-4) are automatically correct.

It would now be possible to rewrite the "algorithm" of Section 6 so that it uses the data structures of 7-1 to 7-7, i.e. using the various reinterpretations set out above.

8 A MODIFICATION

Changes of specification and new ideas to, for example, improve efficiency are facts of life: if Formal Development is to be an accepted way of constructing programs its influence on subsequent changes must be understood. It is the opinion of the current author that, if generality and abstractness have been the keynotes of the development, the difficulty of installing changes should be in realistic relation to how drastic that change is. Certainly more work will be required than with a "quick patch", but the same gains of certainty as with Formal Development should make the investment worthwhile.

The modification made in this section is an optimisation to reduce the number of irrelevant states created and scanned. It is, again, taken from a conventionally built version of the program. The concept is to append to rules an inevitable terminal character (i.e. a character with which any derivation of the rule must start), if such exists. This is specified formally in 8-2. During prediction this look ahead character, if it exists, is compared to the next character of the string: if they are unequal no state for this rule need be generated. Any such state would, after some number of further predictions, have been shown to be a "blind alley" by the scanner. This test will be performed by the predicate "cond" whose essential property is stated in 8-3: an, obvious, realisation is given in 9-13.

The effect of this change is now traced beginning with the definition of is-rule (3-2) which is extended to include the pre-computed look ahead symbols:

(8-1)
$$\text{is-rule} = (\langle \text{s-lhs} : \text{is-nt} \rangle, \\ \langle \text{s-rhs} : \text{is-el-list} \rangle, \\ \langle \text{s-lk} : \text{is-t} \rangle)$$

Further, the restriction on the look ahead character is:

(8-2)
$$\text{for } r \in G : \text{s-lk}(r) \neq \epsilon \supset (\text{s-rhs}(r) \stackrel{*}{\Rightarrow} \alpha \wedge \text{is-t}(\alpha_1) \supset \alpha_1 = \text{s-lk}(r))$$

Now a check of the development process is made to determine the extent of the changes.

Consider the problem definition (Section 3) - there is no change to be made since the whole point is to recognise the same strings, but to do it more efficiently.

Consider the outline of Earley's state sets (Section 4) - there will be no change since the cut down state sets will still fall between those satisfying equations 4-7 and 4-8.

Consider the creation of State sets (Section 5) - A predicate "cond" is introduced which is assumed to satisfy the property:

(8-3)
$$\text{for } r \in G : \text{s-rhs}(r) \stackrel{*}{\Rightarrow} \alpha \wedge \text{is-t}(\alpha_1) \supset \text{cond}(r, \alpha_1)$$

$$\text{cond: RULE} \times T \rightarrow \{T, F\}$$

Then 5-8 is modified to employ cond:

(8-4)
$$\text{predict}(i, X, G, \langle r, j, f \rangle) = \{ \langle s, 0, i \rangle \mid j \neq \bar{r} \wedge \text{is-nt}(\text{s-rhs}(r)_{j+1}) \wedge s \in G \wedge \text{s-lhs}(s) = \\ \text{s-rhs}(r)_{j+1} \wedge \text{cond}(s, X_{i+1}) \}$$

The proof that the state-sets still satisfy 4-6 to 4-9 is modified as follows:

Theorem 5-11	unchanged
Theorem 5-12	unaffected, because the set satisfying 8-4 is clearly a subset of that satisfying 5-8.
Theorem 5-13	slight change related to cond

Consider the mapping of state sets to lists (Section 6) - no change to be made.

Consider the structuring for store (Section 7) - the reinterpretation of 8-1 (see 7-17) includes:

(8-5)
$$\text{s-lk}(r) \quad \text{becomes} \quad \text{s-STRCHAR}(\xi)_{\text{s-N}}(I(\text{s-rhs}(r)_0))$$

and the additional constraint on INGR (see 7-21) -

(8-6)
$$\text{s-lk}(r) = \text{s-STRCHAR}(\xi)_{\text{s-N}}(I(\text{s-rhs}(r)_0))$$

Notice that this simple change to the input properties would require a fairly drastic change to INGR in that space must now be left for all the look ahead characters in s-RULES.

9 CODING IN PL/I

This section discusses the PL/I version (see APPENDIX 1) of the algorithm contained in Section 7 (plus 5-8,5-9,5-10) as modified by Section 8. The transition is not shown in full detail since this would require a significant portion of the theory of PL/I to be formalised. This is not only outside the scope of the current paper, it would also be inappropriate since it should be part of a larger, more general, endeavour.

First a number of constructs of PL/I are related to the properties that have been used above; then a number of specific results are established.

PL/I data properties -

- (9-1) Objects (as in 2-6) can be replaced by PL/I structures and the named selectors can be mimicked by qualified naming. Notice that this property only holds with selectors which cannot yield the null object (see ref 2).
- (9-2) Where the selectors are arbitrary (i.e. their names are not explicitly given) the objects can be replaced by a collection of allocations of a based variable. In this case the selectors are replaced by pointers. (N.B. in the representation of STATE-SETS given, the REFER option has also been used to obtain the possibility to dynamically vary the maximum size of a state set).
- (9-3) Lists (as in 2-10, 2-11) can be replaced by PL/I arrays, these present upper bound problems which make the algorithm non-total. Referencing is via subscripting and the current length is usually stored in another variable.
- (9-4) As an alternative, lists of characters can be replaced by PL/I variable character strings. In this case selection is via the 'SUBSTR' built-in function, and length can be found via the LENGTH built-in function.
- (9-5) Integers are represented by FIXED BINARY numbers of length 15. This again introduces a, difficult to specify, restriction on the domain of the algorithm.

PL/I statement properties -

- (9-6) Conditional Expressions can be replaced by PL/I statements of the form -

```
      IF  p1      THEN ...
      :
      :
      ELSE IF pn  THEN ...
```

- (9-7) With objects which extend monotonically (e.g. state sets) assignment can be used with no loss of information.
- (9-8) The use of where clauses can be replaced by computing the required value and assigning it to a variable which is referred to wherever the name is used (this can be thought of as "assignment as an abbreviation").
- (9-9) PL/I DO loops are used to specify the sequencing of operations. In particular nested DO loops iterating through the collection of state sets and through the members of a state set are used to mirror the function of "next" (see 7-33). Notice that the correctness of this relies on the fact that PL/I re-evaluates the while clause of a DO loop at each iteration. Thus a DO BY TO construct would fail!
- (9-10) The conditional expression form of and (see 2-2) is translated so that the second operand is not evaluated if the first is false.

Special Results -

- (9-11) Lemma $ss' = ss \cup \{<s>\}$ yields same result as a procedure whose argument is a state and which changes the state-set by side effect.
 DO CV = 1 BY 1 TO $\ell(ss)$;
 IF $ss(cv) = <s>$ THEN RETURN ;
 END
 $ss_{\ell(ss)+1} = <s>$;
- (9-12) Lemma The effect of adding a set defined implicitly can be obtained by an ordered search, calling the procedure of 9-10 each time a candidate for addition is located.

$$\text{cond } (r,t) \equiv s\text{-lk}(r) = \text{f} \vee t = s\text{-lk}(r)$$

satisfies 8-3

proof follows from 8-2 and 9-12.

The program invokes procedures called INGR and INSTR to read grammar and string into store. (These routines were in fact taken from the conventionally built program.) A procedure IPTEST is then invoked to check restrictions 4-1,6-9,7-20 and 7-21,8-2. Clearly the relation to the input grammar cannot be checked directly and the expedient of printing the grammar as it "must have looked" is adopted. Finally, the parsing proper is performed.

10 DISCUSSION

This section attempts to discuss some of the open questions. It is the hope that work on Formal Development will be applied to a variety of examples in an attempt to resolve these and other questions.

Breaking the development down into steps has given a clearer picture of the algorithm and its correctness: is the work justified? Making a comparison to the more widely used "write then attempt proof" there are few portions of the above work which are wasted. (One example is the necessity to set up the same inductive form for proofs at two or more steps.) In view of the difficulty of composing a complete algorithm correctly from the beginning, this cost is not great. The proofs also tend to follow the intuitive program ideas more closely. Comparisons to methods not involving formal proofs are more difficult, but the opinion of the current author is that without the discipline of proof such methods will migrate poorly from their original environment.

The comment was made, by an experienced programmer during a presentation of this work, that the much maligned tendency to start coding too early may be a symptom of wanting to write something formal but lacking any alternative language in which to state the partially thought out ideas!

The development given here was aided by the existence and knowledge of an actual working version of the program. In an actual development much more backtracking would be required: it is certainly not claimed that the formal development approach would give rise to inspirations like Earley's state sets. However, the formal approach would certainly offer a framework in which to search - only further work can resolve this point.

One of the principal problems is the level of formality required in the proofs. The developed algorithm was keypunched and run on an Algol grammar. Although no errors have been uncovered this does not establish that the above level of proof would be formal enough to provide security against errors. There is certainly an important distinction to be made between assumptions on the data which are extremely dangerous (being the source of many conventional programming errors), and formality in the deductions which is not attempted. In particular, the extremely attractive, notion of machine checkable proofs may not be practical because of the tedium of formalising the deductive steps. A particular problem of this distinction is the failure in the development to separate variables of the problem and language (e.g. control variables) domains.

There are many general problems which become apparent in the above development: is there a clear set of ways of using storage/assignment (see 9-7,9-8)? Are there better ways of specifying order than the usual DO loops? Questions like this could provide valuable feedback to language design from a problem, rather than the usual machine, angle.

ACKNOWLEDGEMENTS

The author gratefully acknowledges that the idea of levels of an algorithm was suggested by P Lucas. The current paper represents a reasonably complex example and the author's own view of how to add the proofs. Earley's paper provided not only the ideas of the algorithm but also the basis for some of the proofs of the first sections. The author is also grateful for useful discussions with C D Allen, D Chapman, and P D Wright for providing the original running program. Although different in important respects the author acknowledges the stimulus of having read the works of Dijkstra, Hoare and Mills relating to development of programs.

REFERENCES

1. J EARLEY. "An Efficient Context-free parsing Algorithm". Comm ACM Vol 13 No 2, Feb 1970
2. P LUCAS and K WALK. "On the formal Description of PL/I" Annual Review in Automatic Programming Vol 6 Part 3, 1969
3. D PARK. "Fixpoint Induction and Proofs of Program Properties". Machine Intelligence 5, 1969
4. Z MANNA and J McCARTY. "Properties of Programs and Partial Function Logic" Stanford AI Memo 100

```

EARLY:  PROC          OPTIONS(MAIN);
DCL NULL BUILTIN;
      DCL  INPUT CHAR(5000) VARYING INIT('') EXTERNAL;
      DCL  1 NONT(500) EXTERNAL PTR,
        1 STR(250) EXTERNAL,
          2 START FIXED BIN(15),
      STRCHAR CHAR(4000) VARYING INIT('') EXTERNAL,
      1 RULES(4000) EXTERNAL,
        2 TYPE CHAR(1),
        2 N FIXED BIN(15),
      1 RULED BASED(Q),
        2 MAXNRULES FIXED BIN(15),
        2 NEXT      FIXED BIN(15),
        2 RULE(NRULES REFER (MAXNRULES)),
          3 START FIXED BIN(15),
          3 LENGT  FIXED BIN(15);

      DCL IPGR  ENTRY EXTERNAL;
      DCL IPSTR ENTRY EXTERNAL;
      DCL IPTEST ENTRY EXTERNAL;

CALL IPGR  ;
CALL IPSTR ;
CALL IPTEST ;

/* N.B. IPSTR HAS SQUEEZED ALL BLANKS TO AVOID REQUIREMENT TO CODE
   SAME IN ALGOL GRAMMAR */

BEGIN;

DCL S(0:LENGTH(INPUT)) PTR;
DCL 1 STATE BASED (SI),
  2 SIZE FIXED BIN(15),
  2 N      FIXED BIN(15),
  2 INFO (STATE_SIZE REFER (SIZE)),
    3 RULE,
      4 RPTR PTR,
      4 SSC FIXED BIN(15),
    3 RULPOS FIXED BIN(15),
    3 STRPOS FIXED BIN(15);
DCL STATE_SIZE FIXED BIN(15) INIT(70);

DCL I FIXED BIN(15);      /* STATE SET INDEX ,MAJOR LOOP          */
DCL N FIXED BIN(15);      /* STATE INDEX      ,2ND LOOP          */
DCL P1 PTR;               /* PTR OF S-RULE OF STATE(N,STATE-SET(I)) */
DCL I1 FIXED BIN(15);     /* INT OF S-RULE OF STATE(N,STATE-SET(I)) */
DCL J  FIXED BIN(15);     /* RULPOS           OF STATE(N,STATE-SET(I)) */
DCL F  FIXED BIN(15);     /* STRPOS           OF STATE(N,STATE-SET(I)) */
DCL R_ FIXED BIN(15);     /* LEN-PHS(P1,I1)          */

DCL 1 RHS_EL,             /* RHS_EL(J+1,(P1,I1))          */
  2 TYPE CHAR(1),
  2 N FIXED BIN(15);

DCL P2 PTR;               /* PTR PT OF RULES S.T.          */
                           /* RHS_NT EQ LHS((P1,I1),J+1,(P2,?)) */
DCL K FIXED BIN(15);      /* PREDICTOR'S RULE INDEX          */

DCL LK CHAR(1);           /* LK_EL(P2,K)                    */

DCL KS FIXED BIN(15);     /* STATE INDEX WITHIN S(F)          */

DCL P3 PTR;               /* PTR OF S-RULE (STATE(KS,STATE-SET(F)) */
DCL I2 FIXED BIN(15);     /* INT OF S-RULE (STATE(KS,STATE-SET(F)) */
DCL L  FIXED BIN(15);     /* S-RULPOS (STATE(KS,STATE-SET(F)) */
DCL G  FIXED BIN(15);     /* S-STRPOS (STATE(KS,STATE-SET(F)) */

DCL 1 SF_RHS_EL,         /* RHS_EL(P3,I2)                  */
  2 TYPE CHAR(1),
  2 N FIXED BIN(15);

```

```

DCL ADD_STATE ENTRY(PTR, FIXED BIN(15), FIXED BIN(15), FIXED BIN(15),
                    FIXED BIN(15));

DCL SI PTR;          /* CONTAINS S(I) : F COMP RESTR */
DCL SF PTR;          /* CONTAINS S(F) : F COMP RESTR */
DCL NONT1PTR PTR;    /* CONTAINS NONT(1).PTR : F COMP RESTR */

ADD STATE:
PROC (PTR, INT, J, F, I);
    /* ADDS ((PTR, INT), J, F) TO STATE-SET(I) , UNLESS THERE */

DCL PTR PTR;
DCL INT FIXED BIN(15);
DCL J    FIXED BIN(15);
DCL F    FIXED BIN(15);
DCL I    FIXED BIN(15);

DCL K FIXED BIN(15); /* STATE INDEX WITHIN STATE_SET(I) */

DCL SI PTR;          /* CONTAINS S(I) : F COMP RESTR */

SI = S(I);
DO K = 1B BY 1B TO SI -> STATE.N ;
    IF (SI -> STATE.INFO(K).RULE.RPTR = PTR )
        &(SI -> STATE.INFO(K).RULE.SSC = INT )
        &(SI -> STATE.INFO(K).RULPOS = J )
        &(SI -> STATE.INFO(K).STRPOS = F ) THEN
            RETURN; /* STATE ALREADY THERE */
END;

IF SI -> STATE.N = SI -> STATE.SIZE THEN
    DO; /* STATE SET FULL */
        PUT EDIT(' R4 FAILED ')(SKIP, A);
        STOP;
    END;

    ELSE /* ADD */
    DO;
        SI -> STATE.N = SI -> STATE.N + 1B;
        SI -> STATE.INFO(SI -> STATE.N).RULE.RPTR = PTR;
        SI -> STATE.INFO(SI -> STATE.N).RULE.SSC = INT;
        SI -> STATE.INFO(SI -> STATE.N).RULPOS = J;
        SI -> STATE.INFO(SI -> STATE.N).STRPOS = F;
    END;
END;

DO I = 0B BY 1B WHILE(I <= LENGTH(INPUT));
    ALLOCATE STATE; /* SETS SI */
    S(I) = SI;
    SI -> STATE.N = 0B;
END;

SI = S(0B);
SI -> STATE.N = 1B;
SI -> STATE.INFO(1B).RULE.RPTR = NONT(1B);
SI -> STATE.INFO(1B).RULE.SSC = 1B;
SI -> STATE.INFO(1B).RULPOS = 0B;
SI -> STATE.INFO(1B).STRPOS = 0B;

/* INITIAL S S.T. : IS-S(S)
L(STATE SET(0))=1
STATE(1, STATE SET(0)) =
    <(R)(LAS(R)=ROOT(G)), 0, 0>
1 <= I <= L(INPUT) IMP
L(STATE_SET(I, S))= 0 */

```

```

DO I = 0B BY 1B WHILE(I<= LENGTH(INPUT));
SI = S(I);
DO N = 1B BY 1B WHILE(N<= SI -> STATE.N);

P1 = SI -> STATE.INFO(N).RULE.RPTR;
I1 = SI -> STATE.INFO(N).RULE.SSC;
J = SI -> STATE.INFO(N).RULPOS;
F = SI -> STATE.INFO(N).STRPOS;
R_ = P1 -> RULED.RULE(I1).LENGT;

IF(J  $\neq$  R_) &(I  $\neq$  LENGTH(INPUT)) THEN
DO;
  RHS EL = RULES(P1 -> RULED.RULE(I1).START + J ) ;
  IF RHS_EL.TYPE = 'N' THEN
    DO;
      P2 = NONT(RHS_EL.N);
      DO K = 1B BY 1B TO P2 -> RULED.NEXT -1B;
      LK = SUBSTR(STRCHAR,STR(RULES(P2->RULED.RULE(K).START-1B).N)
        .START,1);
      IF(LK=' ') | (LK=SUBSTR(INPUT,I+1B,1)) THEN
        CALL ADD_STATE(P2,K,0B,I,I);
      END;
    END;
  ELSE IF(RHS_EL.TYPE = 'T') THEN
    IF SUBSTR(STRCHAR,STR(RHS_EL.N).START,1) = SUBSTR(INPUT,I+1,1) THEN
      CALL ADD_STATE(P1,I1,J+1B,F,I+1B);
    END;
  ELSE IF J = R_ THEN
    DO;
      SF = S(F);
      DO KS = 1B BY 1B TO SF -> STATE.N;
      P3 = SF -> STATE.INFO(KS).RULE.RPTR;
      I2 = SF -> STATE.INFO(KS).RULE.SSC;
      L = SF -> STATE.INFO(KS).RULPOS;
      G = SF -> STATE.INFO(KS).STRPOS;

      IF L  $\neq$  P3->RULED.RULE(I2).LENGT THEN
        DO;
          SF_RHS_EL = RULES(P3->RULED.RULE(I2).START+L);
          IF SF_RHS_EL.TYPE = 'N' THEN
            IF NONT(SF_RHS_EL.N) = P1 THEN
              CALL ADD_STATE(P3,I2,L+1B,G,I);
            END;
          END;
        END;
      END;
    END;
  END;
END;
END;
END;

/* CHECK FOR END_STATE IN S(L(X)) */

SI = S(LENGTH(INPUT));
NONT1PTR = NONT(1B);

DO N = 1B BY 1B TO SI -> STATE.N;
IF (SI -> STATE.INFO(N).RULE.RPTR = NONT(1B))
&(SI -> STATE.INFO(N).RULE.SSC = 1B )
&(SI -> STATE.INFO(N).RULPOS = NONT1PTR -> RULED.RULE(1B).LENGT)
&(SI -> STATE.INFO(N).STRPOS = 0B ) THEN
DO;
  PUT EDIT (' YES' ) (SKIP ,A); RETURN ;
END;
END;

PUT EDIT (' NO' ) (SKIP,A); RETURN ;
END;
END;

```