

*Discussed is the "syntax machine," a program for automatically generating syntactically correct programs (test cases) for checking compiler front ends.*

*The notion of "dynamic grammar" is introduced and is used in a syntax-defining notation that provides for context-sensitivity.*

*Examples demonstrate use of the syntax machine.*

## **Automatic generation of test cases**

**by K. V. Hanford**

The "syntax machine" discussed here automatically generates random test cases for any suitably defined programming language.<sup>1</sup> The test cases it produces are syntactically valid programs. But they are not "meaningful," and if an attempt is made to execute them, the results are unpredictable and uncheckable. For this reason, they are less valuable than handwritten test cases. However, as an inexhaustible source of new test material, the syntax machine has shown itself to be a valuable tool.

In the following sections, we characterize the use of this tool in testing different types of language processors, introduce the concept of "dynamic grammar" of a programming language, outline the structure of the system, and show what the syntax machine does by means of some examples.

### **Test cases**

Test cases for a language processor are programs written following the rules of the language, as documented. The test cases, when processed, should give known results. If this does not happen, then either the processor or its documentation is in error.

We can distinguish three categories of language processors and assess the usefulness of the syntax machine for testing them. For an interpreter, the syntax machine test cases are virtually useless,

since the result of their execution is unpredictable. For a syntax checker, the syntax machine can do a complete job, including the testing of the diagnostic facilities. For a compiler, the syntax machine can fully test the ability of the input phase to accept syntactically valid input, but does not test the translation process. It can be used to estimate the reliability of the compiler with respect to such problems as getting into an infinite loop, abnormally ending, or diagnosing nonexistent syntax errors.

Although as a writer of test cases, the syntax machine is certainly unintelligent, it is also uninhibited. It can test a processor with many combinations that would not be thought of by a human test case writer. Also, test cases can be created at will, so that at any time new material can be presented to a product under test. This is particularly important for a product that has "learned" its test case library.

In its standard mode, the syntax machine produces programs in a pseudo-random manner. The system has a second mode in which it systematically produces all the possible syntactic forms of certain designated language features. The system can also be directed to produce programs that contain syntactic errors, with annotations giving the position of the errors.

### Dynamic grammar

This section discusses how declarations can be considered to modify syntax and how formal notation can be extended to account for this.

Consider a compiler for an algebraic programming language. Its job is to manufacture sequences of elementary machine instructions to simulate the evaluation of expressions in the language. To such a compiler,  $2 \times (3 + 4)$  is acceptable for translation—the required simulating sequence can be determined without difficulty. On the other hand,  $2 \times 3 + 4$  and  $2 \times 3 (+ 4$  are not acceptable for translation, since the compiler cannot confidently ascribe a meaning to them. Thus we say that the first expression is "well-formed" or "grammatically correct," whereas the other two are not.

Let us define the "acceptable programs" for a given compiler to be those source programs whose compilation is trouble-free and results in no diagnostic messages. Then we find that acceptance depends on form rather than meaning: a compiler would accept

$$x = 0$$

$$y = 1/x$$

on its form, although it is computationally meaningless. Thus, the notion of "well-formed" applies not only to algebraic expres-

sions but to all the elements of a programming language, including whole statements and whole programs. If we compare two compilers for the same language, we find that the set of acceptable programs for one is almost identical to the set of acceptable programs for the other. There are usually some gray areas of difference, but these serve only to emphasize the high level of agreement on what is acceptable and what is not.

This leads us to ask whether the concept of grammar in programming language can be formalized. Can we make a definition of "well-formed" or "grammatically correct" that is a good, useful approximation to the pragmatic notion of "acceptable" programs or, alternatively, that could be used as an exact official specification of those source programs that should be accepted without diagnostic messages?

The first formal definition of the syntax of a programming language was made by Backus<sup>2</sup> in a notation now known as BNF (Backus normal or Backus-Naur form). A class of strings of a language is given a name enclosed in brackets  $\langle \rangle$ , and a class is defined by one or more "writing rules." An example of a rule is:

$$\langle x \rangle \rightarrow \langle y \rangle + (\langle z \rangle)$$

This says that if you write down *any* member of the class  $\langle y \rangle$ , followed by the mark +, followed by the mark (, followed by *any* member of the class  $\langle z \rangle$ , followed by the mark ), then what you have written down is some member of the class  $\langle x \rangle$ . The classes  $\langle y \rangle$  and  $\langle z \rangle$  must in their turn be defined by further writing rules. The syntax of a language consists of a finite number of such rules.

Following the work of Backus (related to earlier work by Chomsky<sup>3</sup>), a formal approach to syntax in the definition of programming languages has become commonplace. Notations vary but the motivation is always the same—to define as fully as possible the form of a program.

The grammars found in language definitions are recipes for writing down strings. The intention is that the recipe should yield precisely the set of strings that will be considered to be well-formed and must be accepted by a compiler. To meet this intention fully, the following two conditions need to be satisfied:

1. *All* well-formed programs can be written down following the recipe.
2. *Only* well-formed programs can be written down following the recipe.

BNF meets the first condition but fails to meet the second. The set of strings that it defines certainly contains all well-formed programs

but contains other strings as well. In language definition documents, to further restrict the set of defined strings, the BNF grammar is accompanied by extra grammatical rules, stated in English. These rules consist of statements like: "The constituents of an arithmetic expression must have been declared to be of type real or type integer." These extra rules are given in English because it is impossible to express them in BNF.

The informal grammatical rule given above says that a variable may occur in an arithmetic expression if and only if it also occurs in some program element declaring it to be of type real or type integer. Thus, what can be written at one place in a program depends in general on what has been written down in other places. A grammatical rule of this kind, which requires relations to hold between parts of a program, is called a "context sensitivity," and languages that have rules of this kind are "context-sensitive." All present-day high-level programming languages are context-sensitive.

The set of strings defined by a BNF grammar is termed a "context-free" language. A BNF formulation for a high-level language defines a context-free language that is a superset of the actual (context-sensitive) language.

The context-sensitive nature of programming languages derives from the fact that these languages allow *declarations*. The programmer uses a declaration to create an object, to give it a name, and to describe what kind of object it is, e.g., an array, a function, etc. Now, the contexts in which an object may occur in a program depend on the kind of object it is. Thus, declarations ascribe certain syntactic qualities to objects and the declarations in a program influence the syntax of that particular program. For example, the real and integer type declarations of a program determine what arithmetic expressions can be written in that program.

The syntax machine involves a method for the description of context-sensitive constraints. The basic idea in our formulation of these constraints is that the effect of writing a declarative statement in a program can be represented as a *change to the context-free grammar*. The rules for writing down a program should, therefore, not only specify what can be written down but also specify what changes to the grammar result from what is written down. By this means, the declarative statements construct the grammar for the other statements. This concept can be stated as follows:

The programs of a declarative programming language can be produced from a context-free grammar that has the ability to modify itself. A grammar of this form will be called a "dynamic grammar."

As an example, consider the FORTRAN declaration

INTEGER XYZ

If this string is written down as a statement of a program, it causes the grammar to be modified by the addition of a rule of the form

$\langle \text{integer variable} \rangle \rightarrow \text{XYZ}$

Initially, the grammar contains no rules for  $\langle \text{integer variable} \rangle$ . Each declaration of the above form causes a new rule for  $\langle \text{integer variable} \rangle$  to be added to the grammar. For example, the declarations

INTEGER AY  
INTEGER BEE  
INTEGER C

cause the dynamic grammar to be enriched by the rules:

$\langle \text{integer variable} \rangle \rightarrow \text{AY}$   
 $\langle \text{integer variable} \rangle \rightarrow \text{BEE}$   
 $\langle \text{integer variable} \rangle \rightarrow \text{C}$

When subsequently an arithmetic expression is written, these rules are available for use.

### Production systems

In this section, we discuss first the requirements of a system for producing the strings of a context-free language. We then consider how this system can be extended to a context-sensitive production system by using syntax generators with a rewriting algorithm. A more detailed description of the rewriting procedure is then presented, which takes into account the contingencies that may be encountered in rewriting a nonterminal.

#### context-free production system

Consider a context-free language defined by a BNF grammar. This grammar consists of a set of rules such as

$\langle x \rangle \rightarrow \langle y \rangle + (\langle z \rangle)$

The left-hand side is a class-name or "nonterminal" and the right-hand side is a sequence of symbols and nonterminals.

A process for obtaining a string of the class  $\langle x \rangle$  is described by the following "rewriting" scheme. Starting with the nonterminal  $\langle x \rangle$ , we rewrite it using one of the rules for  $\langle x \rangle$ . Suppose we choose the above rule. The right-hand side of the rule is a recipe. It uses two methods for obtaining a sequence of symbols to form a string of the  $\langle x \rangle$  class. It directly contributes some symbols (the symbols  $+$ ,  $($ ,  $)$  in the example). It obtains other symbols indirectly by rewriting the nonterminals that occur in the recipe ( $\langle y \rangle$ ,  $\langle z \rangle$  in the example). In general, there is a choice of recipes for the rewriting of a given nonterminal. When a subsidiary recipe has done its job, it returns responsibility for symbol production back to the recipe that called it. To obtain a complete string of the

language, we start the process with the "language nonterminal" (e.g.,  $\langle \text{program} \rangle$ ), which names the class of complete strings.

Thus, the process of writing down a program of a context-free programming language can be mechanized as a rewriting algorithm that interprets a BNF grammar. A rule of the grammar is scanned from left to right. When a nonterminal is met, a pseudo-random choice is made from the rules for that nonterminal and a recursive call to the rewriting algorithm is made. The context-free grammar and the rewriting algorithm together constitute a "production system" for a context-free language.

A "dynamic production system" is a generalization of a context-free production system and again consists of a set of production rules and a rewriting algorithm. The generalizations concern:

**context-sensitive  
production  
system**

- The "activation" during the production process of "syntax generators" for the self-modification of the production system
- A cyclic rewriting algorithm using "delayed nonterminals"

The rewriting algorithm is an extension of the rewriting algorithm for a context-free production system. As we have seen, context-sensitivity constraints are considered to arise from declarations. The effect of a declaration is represented by a change to the context-free grammar. When a declaration is produced, the context-free grammar is immediately modified. This modification is carried out in the following fashion. Certain context-free production rules, generally those concerned with declarations, have syntax generators attached to them. The production algorithm, after using a production rule for rewriting, activates any syntax generator attached to that rule. Syntax generators synthesize new rules and add them to the context-free grammar.

Now a declaration may influence that part of the program that comes before it. For example, we may have a branch statement with a forward reference to a label (the defining occurrence of a label is considered to be an implicit declaration). This problem of "use before declaration" is handled by giving the "delay" qualification to any nonterminal whose rewriting may depend on the rewritten form of nonterminals occurring to its right. When a nonterminal with the delay qualification is met in the left-to-right rewriting scheme, its delay qualification is removed, but it is not rewritten until the remaining nonterminals to its right have been rewritten. Delayed nonterminals are indicated by the character  $-$ ; thus  $\langle -x \rangle$  stands for the nonterminal  $\langle x \rangle$  with the delay qualification.

In the following more detailed description of the rewriting algorithm, we discuss the conditions under which syntax generators are activated but not how they create rules. This topic is illustrated by

example in the next section. "Not rules" mentioned in the description are of the same form as rules and mean that a string of the class named by the left-hand nonterminal "cannot be" a member of the class defined by the right-hand expression.

The description of the rewriting algorithm is divided into discussions of the recursive rewriting of a nonterminal, backtracking, and the cyclic rewriting of a complete program.

**recursive  
rewriting**

Let  $\langle x \rangle$  be a nonterminal element of an arbitrary string. The recursive rewriting of  $\langle x \rangle$  ought to replace  $\langle x \rangle$  by a string. However, several factors can affect this effort to rewrite  $\langle x \rangle$ , including the presence of the delay qualification mentioned above, the possible existence of "not rules," the absence of any rules for  $\langle x \rangle$ , and the possibility that a nonterminal in a rule chosen for  $\langle x \rangle$  can itself not be rewritten.

We suppose that  $\langle x \rangle$  is not a delayed nonterminal. A rule for  $\langle x \rangle$  is chosen randomly, and  $\langle x \rangle$  is replaced by the right-hand side of this rule with all its nonterminals recursively rewritten in left-to-right sequence.

If there are no not rules for  $\langle x \rangle$  and if the chosen rule for  $\langle x \rangle$  has a syntax generator, the generator is activated. This completes the recursive rewriting of  $\langle x \rangle$ .

Now suppose there exist not rules for  $\langle x \rangle$ . Then the string that has replaced  $\langle x \rangle$  is tested for membership of the "not class" for  $\langle x \rangle$ . If the string is a member of the not class, it is rejected (and the generator is not activated). The algorithm tries again to rewrite  $\langle x \rangle$ , using a new randomly chosen rule. Only when a string is obtained that is not a member of the not class is it accepted. The generator of the successful rule is then activated.

**back-  
tracking**

There are three possible ways in which the process described above can fail to rewrite a given nonterminal  $\langle x \rangle$ . This possibility of failure leads to the need for a "backtracking" facility in the rewriting algorithm. These three possible ways are:

1. There are no rules for  $\langle x \rangle$ .
2. The strings obtained for  $\langle x \rangle$  are consistently members of the not class for  $\langle x \rangle$ .
3. Some nonterminal in the chosen rule for  $\langle x \rangle$  fails to be rewritten.

In order to explain what happens in these situations, let us suppose that the writing of  $\langle x \rangle$  was called for (directly) in the rewriting of  $\langle y \rangle$ . That is, the chosen rule for  $\langle y \rangle$  is of the form

$$\langle y \rangle \rightarrow \dots \langle x \rangle \dots$$

In Case 1, the production algorithm backtracks. It abandons the attempt to write  $\langle x \rangle$  and goes back to the  $\langle y \rangle$  level. With respect to this level, Case 3 holds.

In Case 2, 25 attempts are made to find a string that is not a member of the not class of  $\langle x \rangle$ . At each attempt, a new random choice of a rule is made. After 25 unsuccessful attempts, the attempt to rewrite  $\langle x \rangle$  is abandoned and the algorithm behaves as if there are no rules for  $\langle x \rangle$  (Case 1).

In Case 3, the production algorithm makes a new random choice of a rewriting rule for  $\langle x \rangle$ , but debars the abortive rule from being chosen. If Case 3 also arises for the new rule, it in turn is debarred from being chosen at the next attempt. If all the rules for  $\langle x \rangle$  are eventually debarred in this way, the attempt to rewrite  $\langle x \rangle$  is abandoned and the algorithm behaves as if there are no rules for  $\langle x \rangle$  (Case 1).

There are two situations in which backtracking can give erroneous results. First, suppose that backtracking goes back to a nonterminal  $\langle z \rangle$ , say, whose rewriting involved the activation of a generator, e.g., the chosen rule for  $\langle z \rangle$  was of the form

$$\langle z \rangle \rightarrow \dots \langle x \rangle \dots \langle y \rangle \dots$$

where the chosen rule for  $\langle x \rangle$  had a generator and  $\langle y \rangle$  could not be written. The effect of the generator is irreversible, so that an invalid change has probably occurred in the grammar. In this situation, a diagnostic message is given, but the production continues. Second, suppose the production process fails to rewrite a nonterminal that was delayed in the previous pass. Then backtracking will not occur and the attempted production is terminated; the syntax machine is unable to backtrack past the start of its current rewriting pass.

Let  $\langle x \rangle$  be a nonterminal. The "cyclic rewriting" of  $\langle x \rangle$  replaces  $\langle x \rangle$  with a string of symbols, as follows.  $\langle x \rangle$  is recursively rewritten. If the string that replaces  $\langle x \rangle$  consists only of symbols, the process is complete. Otherwise, the remaining nonterminals are recursively rewritten, in left-to-right sequence. This step is iterated until no nonterminals remain.

cyclic  
rewriting

A language string is obtained by cyclically rewriting the language nonterminal.

The syntax machine is such a generalized production system. The actual program consists of two parts: a grammar loader and a production algorithm. The loader stores the grammar in main storage in a form convenient for use by the production algorithm in creating programs.



Table 1 Syntax of Little PL/I

<i>Symbols, identifiers and numbers</i>		
<letter>	→ a b ... z	1.
<digit>	→ 0 1 ... 9	2.
<identifier>	→ <letter>   <identifier> <letter>	3.
<unsigned integer>	→ <digit>   <unsigned integer> <digit>	4.
<i>Fixed-point variables and expressions</i>		
<fixed-point variable>	→ i j k l m n  <fixed-point variable> <letter>	5.
<primary>	→ <unsigned integer>   <fixed-point variable>   ( <arithmetic expression> )	6.
<relational operator>	→ <  =  >	7.
<boolean expression>	→ <primary> <relational operator> <primary>	8.
<arithmetic operator>	→ + - * /	9.
<arithmetic expression>	→ <primary>   <arithmetic expression> <arithmetic operator> <primary>	10.
<i>Labels and label variables</i>		
<label>	→ <identifier>	11.
<label variable>	→ <identifier>	12.
<i>Statements</i>		
<arithmetic assign statement>	→ <fixed-point variable> = <arithmetic expression>;	13.
<label assign statement>	→ <label variable> = <label>;   <label variable> = <label variable>;	14.
<go to statement>	→ go to <label>;   go to <label variable>;	15.
<if statement>	→ if <boolean expression> then <nondeclare statement>	16.
<nondeclare statement>	→ <arithmetic assign statement>   <label assign statement>   <go to statement>   <if statement>	17.
<declare statement>	→ declare <label variable> label;	18.
<statement>	→ <nondeclare statement>   <label> : <nondeclare statement>   <declare statement>   <label> : <declare statement>	19.
<i>Programs</i>		
<statement sequence>	→ <statement>   <statement sequence> <statement>	20.
<pl/i program>	→ <label> : procedure; <statement sequence> end <label>;	21.

### Using the system

#### notation

Table 2 Little PL/I program due to Donovan and Ledgard

```

Q: PROCEDURE;
  DECLARE LX LABEL;
  L: I=I+IA×IB-IC;
    LX=L;
    GO TO CHECK;
  M: I=I+1;
    LX=M;
  CHECK: IF I< LIMIT
    THEN GO TO LX;
END Q;

```

To illustrate the use of the syntax machine, we now develop a definition for the syntax of a subset of PL/I. This subset has been taken from a paper by Donovan and Ledgard,<sup>4</sup> who call it Little PL/I. It includes limited forms of the following types of statements: GOTO, IF, label declaration, label assignment, and arithmetic assignment. A BNF definition for the context-free syntax of this language is given in Table 1. The mark | is used to separate the alternatives for a given class.

The following quotation is taken from the Donovan and Ledgard paper. The quotation refers to the example of a Little PL/I program shown in Table 2.

"We define the syntax of a language as the set of rules for specifying the strings that can be recognized by a translator and translated

into some other language. The set of rules excludes strings that the translator would reject solely on their form. The syntax of Little PL/I has the following restrictions, which for all practical purposes make Little PL/I context-sensitive and therefore impossible to completely characterize in Backus-Naur Form:

1. Different declarations of the same identifier are in error, i.e.,
  - a. The lists of fix-pt variables, statement labels, and declared label variables for a program must be mutually disjoint;
  - b. The label before PROCEDURE must not occur within the procedure block.
2. The label after END must be identical to the label before PROCEDURE.
3. All statement labels must be different.
4. The identifier in a GOTO statement must refer to an existing statement label or a declared label variable.
5. The identifier on the left side of the = in a label assignment statement must refer to a declared label variable; the identifier on the right side of the = must refer to an existing statement label or a declared label variable."

Figure 1 gives a basic grammar for Little PL/I in the notation of the syntax machine. Upper case characters have been listed as lower case characters.

The input is composed of system statements.<sup>5</sup> Each statement begins on a new card and starts with a "master phrase." The end of the statement is indicated by the appearance of the next master phrase (this requires that the master phrases be distinct from all other character sequences at the start of a card). There are two classes of system statement, syntax statements and control statements. Syntax statements are used in defining a grammar. The only control statements in Figure 1 are job and comment statements, with obvious meaning. Further control statements will be introduced in the course of our development of the Little PL/I definition.

Blanks are ignored. Class names are enclosed in brackets <>. Basic symbols are represented by themselves except for the blank and the bracket <, which are represented by the pseudo class names <blank> and <less than>, respectively. This is necessary because these two marks have a special use, viz., blank is ignored and the bracket < introduces a class name. The notation also contains three nongraphic pseudo class names: <eor> terminates a punched card, <eol> terminates a printed line, and <nil> represents the empty string.

The section numbers of Table 1 and Figure 1 correspond. For most sections, Figure 1 is merely a transliteration of the BNF of Table 1. However, for certain sections, the syntax machine defini-

Figure 1 Basic grammar of Little PL/I

:job:	10	:rule:	<digit> -> 9	2, 1220
:comment:	20	:rule:	<identifier> -> <letter>	3, 1230
	30	:begin:	<<lambda>> -> <<l>>	3, 1240
	40	:onlyrule:	<<lambda>> -> <<l>>	3, 1250
:begin:	50	:end:		3, 1260
:comment:	60	:rule:	<identifier> -> <identifier> <letter>	3, 1270
	70	:begin:	<<lambda>> -> <<lambda>> <<l>>	3, 1280
	80	:rule:	<<lambda>> -> <<lambda>> <<l>>	3, 1290
	90	:end:		3, 1300
:rule:	1, 100	:rule:	<declaration identifier> -> <identifier>	3, 1310
:begin:	1, 110	:begin:	<declared identifier> -> <lambda>>	3, 1320
:onlyrule:	1, 120	:rule:	<declared identifier> -> <lambda>>	3, 1330
:end:	1, 130	:end:		3, 1340
:rule:	1, 140	:notrule:	<declaration identifier> -> <declared identifier>	3, 1350
:begin:	1, 150	:rule:	<unsigned integer> -> <digit>	3, 1360
:onlyrule:	1, 160	:rule:	<unsigned integer> -> <unsigned integer> <digit>	3, 1370
:end:	1, 170	:comment:	fixed-point variables and expressions	3, 1380
:rule:	1, 180		*****	3, 1390
:begin:	1, 190	:rule:	<i_n identifier> -> i	4, 1400
:onlyrule:	1, 200	:rule:	<i_n identifier> -> j	4, 1410
:end:	1, 210	:rule:	<i_n identifier> -> k	4, 1420
:rule:	1, 220	:rule:	<i_n identifier> -> l	4, 1430
:begin:	1, 230	:rule:	<i_n identifier> -> m	4, 1440
:onlyrule:	1, 240	:rule:	<i_n identifier> -> n	4, 1450
:end:	1, 250	:rule:	<i_n identifier> -> <i_n identifier> <letter>	4, 1460
:rule:	1, 260	:rule:	<fixed_point variable> -> <i_n identifier>	4, 1470
:begin:	1, 270	:notrule:	<fixed_point variable> -> <declared identifier>	4, 1480
:onlyrule:	1, 280	:rule:	<primary> -> <unsigned integer>	4, 1490
:end:	1, 290	:rule:	<primary> -> <fixed_point variable>	4, 1500
:rule:	1, 300	:rule:	<primary> -> ( <arithmetic expression> )	4, 1510
:begin:	1, 310	:rule:	<relational operator> -> <less than>	4, 1520
:onlyrule:	1, 320	:rule:	<relational operator> -> =	4, 1530
:end:	1, 330	:rule:	<relational operator> -> >	4, 1540
:rule:	1, 340	:rule:	<boolean expression> -> <primary>	4, 1550
:begin:	1, 350	:rule:	<boolean expression> -> <relational operator>	4, 1560
:onlyrule:	1, 360	:rule:	<boolean expression> -> <primary>	4, 1570
:end:	1, 370	:rule:	<arithmetic operator> -> +	4, 1580
:rule:	1, 380	:rule:	<arithmetic operator> -> -	4, 1590
:begin:	1, 390	:rule:	<arithmetic operator> -> *	4, 1600
:onlyrule:	1, 400	:rule:	<arithmetic operator> -> /	4, 1610
:end:	1, 410	:rule:	<arithmetic expression> -> <primary>	4, 1620
:rule:	1, 420	:rule:	<arithmetic expression> -> <arithmetic expression>	4, 1630
:begin:	1, 430	:rule:	<arithmetic expression> -> <arithmetic operator>	4, 1640
:onlyrule:	1, 440	:rule:	<arithmetic expression> -> <primary>	4, 1650
:end:	1, 450	:comment:	labels and label variables	4, 1660
:rule:	1, 460		*****	4, 1670
:begin:	1, 470	:rule:	<label assign statement> -> <label variable> = <label>	4, 1680
:onlyrule:	1, 480	:rule:	<label assign statement> -> <label variable> =	4, 1690
:end:	1, 490	:rule:	<label variable> =	4, 1700
:rule:	1, 500	:rule:	<go to statement> -> goto <blank> <label>	4, 1710
:begin:	1, 510	:rule:	<go to statement> -> goto <blank> <label variable>	4, 1720
:onlyrule:	1, 520	:rule:	<if statement> -> if <blank> <boolean expression> <blank>	4, 1730
:end:	1, 530	:rule:	then <blank> <non_declare statement>	4, 1740
:rule:	1, 540	:rule:	<non_declare statement> -> <arithmetic assign statement>	4, 1750
:begin:	1, 550	:rule:	<non_declare statement> -> <label assign statement>	4, 1760
:onlyrule:	1, 560	:rule:	<non_declare statement> -> <go to statement>	4, 1770
:end:	1, 570	:rule:	<non_declare statement> -> <if statement>	4, 1780
:rule:	1, 580	:rule:	<label variable declaration> -> <declaration identifier>	4, 1790
:begin:	1, 590	:rule:	<label variable> -> <<lambda>>	4, 1800
:onlyrule:	1, 600	:rule:	<declare statement> -> declare <blank> <label variable>	4, 1810
:end:	1, 610	:rule:	declaration <blank> label ;	4, 1820
:rule:	1, 620	:rule:	<label declaration> -> <declaration identifier>	4, 1830
:begin:	1, 630	:rule:	<label> -> <<lambda>>	4, 1840
:onlyrule:	1, 640	:rule:	<statement> -> <- non_declare statement> <eor> <eor>	4, 1850
:end:	1, 650	:rule:	<statement> -> <label declaration> ;	4, 1860
:rule:	1, 660	:rule:	<statement> -> <- non_declare statement> <eor> <eor>	4, 1870
:begin:	1, 670	:rule:	<statement> -> <declare statement> <eor> <eor>	4, 1880
:onlyrule:	1, 680	:rule:	<statement> -> <declare identifier> ;	4, 1890
:end:	1, 690	:rule:	<statement> -> <declare statement> <eor> <eor>	4, 1900
:rule:	1, 700	:rule:	<statement sequence> -> <statement>	4, 1910
:begin:	1, 710	:rule:	<statement sequence> -> <statement sequence> <statement>	4, 1920
:onlyrule:	1, 720	:rule:	<procedure name declaration> -> <declaration identifier>	4, 1930
:end:	1, 730	:rule:	<procedure name> -> <<lambda>>	4, 1940
:rule:	1, 740	:rule:	<procedure name declaration> -> <procedure name>	4, 1950
:begin:	1, 750	:rule:	<pl/i program> -> <set up> <procedure name declaration>	4, 1960
:onlyrule:	1, 760	:rule:	procedure ; <eor> <eor>	4, 1970
:end:	1, 770	:rule:	<+ statement sequence> end <blank>	4, 1980
:rule:	1, 780	:rule:	<procedure name> ;	4, 1990
:begin:	1, 790	:rule:	<set up> -> <nil>	4, 2000
:onlyrule:	1, 800	:rule:	<declared identifier>	4, 2010
:end:	1, 810	:rule:	<label>	4, 2020
:rule:	1, 820	:rule:	<label variable>	4, 2030
:begin:	1, 830	:rule:	<procedure name>	4, 2040
:onlyrule:	1, 840	:rule:	<procedure name>	4, 2050
:end:	1, 850	:rule:	<procedure name>	4, 2060
:rule:	1, 860	:rule:	<procedure name>	4, 2070
:begin:	1, 870	:rule:	<procedure name>	4, 2080
:onlyrule:	1, 880	:rule:	<procedure name>	4, 2090
:end:	1, 890	:rule:	<procedure name>	4, 2100
:rule:	1, 900	:rule:	<procedure name>	4, 2110
:begin:	1, 910	:rule:	<procedure name>	4, 2120
:onlyrule:	1, 920	:rule:	<procedure name>	4, 2130
:end:	1, 930	:rule:	<procedure name>	4, 2140
:rule:	1, 940	:rule:	<procedure name>	4, 2150
:begin:	1, 950	:rule:	<procedure name>	4, 2160
:onlyrule:	1, 960	:rule:	<procedure name>	4, 2170
:end:	1, 970	:rule:	<procedure name>	4, 2180
:rule:	1, 980	:rule:	<procedure name>	4, 2190
:begin:	1, 990	:rule:	<procedure name>	4, 2200
:onlyrule:	1, 1000	:rule:	<procedure name>	4, 2210
:end:	1, 1010	:rule:	<procedure name>	4, 2220
:rule:	1, 1020	:rule:	<procedure name>	4, 2230
:begin:	1, 1030	:rule:	<procedure name>	4, 2240
:onlyrule:	1, 1040	:rule:	<procedure name>	4, 2250
:end:	1, 1050	:rule:	<procedure name>	4, 2260
:rule:	1, 1060	:rule:	<procedure name>	4, 2270
:begin:	1, 1070	:rule:	<procedure name>	4, 2280
:onlyrule:	1, 1080	:rule:	<procedure name>	4, 2290
:end:	1, 1090	:rule:	<procedure name>	4, 2300
:rule:	1, 1100	:rule:	<procedure name>	4, 2310
:begin:	1, 1110	:rule:	<procedure name>	4, 2320
:onlyrule:	1, 1120	:rule:	<procedure name>	4, 2330
:end:	1, 1130	:rule:	<procedure name>	4, 2340
:rule:	2, 1140	:rule:	<procedure name>	4, 2350
:begin:	2, 1150	:rule:	<procedure name>	4, 2360
:onlyrule:	2, 1160	:rule:	<procedure name>	4, 2370
:end:	2, 1170	:rule:	<procedure name>	4, 2380
:rule:	2, 1180	:rule:	<procedure name>	4, 2390
:begin:	2, 1190	:rule:	<procedure name>	4, 2400
:onlyrule:	2, 1200	:rule:	<procedure name>	4, 2410
:end:	2, 1210	:rule:	<procedure name>	4, 2420

tion is an elaboration of the BNF definition, where the new material concerns the context-sensitive constraints of Little PL/I. Note also that, in Figure 1, sections 11 and 12 are empty—originally there are no rules for <label> and <label variable>.

The concept of dynamic syntax is implemented in the syntax machine by allowing "syntax generators" to be attached to the rules of a BNF grammar. These generators require an "environment" for use as a working store. An environment consists of "metaclasses" whose names are written between double pointed brackets <<, >>. When a rule has a syntax generator attached to it, this is given following the rule and is bracketed by begin and end statements. An example of a generator occurs in section 19:

```
:rule:  <label declaration> → <declaration identifier>
      :begin:
      :rule:  <label> → <<lambda>>
      :end:
```

The first line is a rule of the context-free grammar, and the following lines are the associated syntax generator. When the rule is used and produces the label declaration abc (say), the syntax generator is "activated," causing the new rule <label> → abc to be added to the context-free grammar.

In the grammar for identifiers, sections 1 and 3, the generators are designed to place a copy of a produced identifier in the metaclass <<lambda>>. Consider, for example, the production of the identifier abc. The rules for <identifier> use left recursion, and abc is formed piecemeal as a, ab, abc. This results in the activation of the sequence of generators.

1. :only rule: <<1>> → a
2. :only rule: <<lambda>> → <<1>>
3. :only rule: <<1>> → b
4. :only rule: <<lambda>> → <<lambda>> <<1>>
5. :only rule: <<1>> → c
6. :only rule: <<lambda>> → <<lambda>> <<1>>

An "only rule" for the metaclasses <<1>> or <<lambda>> erases the old values before assigning the new values. At the end of the above sequence of generator activations, <<lambda>> has the single value abc.

The grammar for <declaration identifier> also contains the "not rule":

```
:not rule: <declaration identifier> → <declared identifier>
```

The definition of a class <x> is that it consists of all those strings that are allowed by the rules for <x> and are not disallowed

by the not rules for  $\langle x \rangle$ . (The class of strings defined by the not rules for  $\langle x \rangle$  is called the "not class" for  $\langle x \rangle$ .) By the use of not rules, we can arrange that the same identifier cannot be declared twice.

We saw above how section 19 arranges that the implicit declaration of a label  $\alpha$ , say, causes the rule  $\langle \text{label} \rangle \rightarrow \alpha$  to be added to the grammar. Similarly sections 18 and 21 arrange that the declaration of a label variable and the procedure name cause rules for  $\langle \text{label variable} \rangle$  and  $\langle \text{procedure name} \rangle$  to be added to the grammar. Initially there are no rules for  $\langle \text{label} \rangle$ ,  $\langle \text{label variable} \rangle$ , and  $\langle \text{procedure name} \rangle$ ; rules for these are created as a result of declarations.

In section 99, the "set-up generator" attached to a dummy rule ensures that the dynamic part of the grammar is initialized to be empty at the start of a program production.

### Examples

We now describe the effect of providing to the syntax machine a number of variations on the definition of Little PL/I given in Figure 1. Updates to the definition are given card sequence numbers to show where they are to be inserted.

```
:rewrite:  $\langle + \text{pl/i program} \rangle$  10      4000
```

The system will produce 10 programs. (If the + were omitted, only one rewrite cycle would be performed for each production, so that the output strings would contain some nonterminals that had not been rewritten.)

```
:weight: 1      1685
:weight: 10     1715
```

The first of these statements assigns "weights" of 1 to the second, third, and fourth rules for  $\langle \text{arithmetic operator} \rangle$ . The second statement resets the "weight register" to its default value of 10. The current value of the weight register is used in assigning a weight to each new rule. In the 10 produced programs, the operators +, -, \*, / will occur with relative frequencies of approximately 10/13, 1/13, 1/13, and 1/13, respectively.

```
:rule:  $\langle \text{statement sequence} \rangle \rightarrow \langle \# 10, 20, 25 \rangle \langle \text{statement} \rangle$  2220
      remove this statement                                     2230
```

In the 10 produced program, the procedure body will consist of from 10 to 25 statements, with a mean of 20.

$\langle \# 10, 20, 25 \rangle$  is an "iteration nonterminal"; these may be used

Table 3 Example of output obtained with Little PL/I

```

SCS:PROCEDURE;
DECLARE IAK LABEL;
N:=1155;
IF 2>>(I1) THEN IF (3)>9 THEN N:=03+26+K+K;
ES:DECLARE ZMQZW LABEL;
XN:DECLARE E LABEL;
IF 5>27 THEN J:=2;
DECLARE EG LABEL;
DECLARE J LABEL;
HMKHHEPAVVC:=(K)/27+1;
MKBP:DECLARE JK LABEL;
IF (N+4+I*(1+2+4+K+KJ)-1)>(7*(L/L+((311)/((35+7*((N+4)))))) THEN IF 53> THEN IF 1CN THEN K:=(1)+7;
C:DECLARE IOLTVJ LABEL;
DECLARE T LABEL;
F:DECLARE M LABEL;
GOTO HMKHHEPAVVC;
ITV:DECLARE Z LABEL;
X:IF (11)>1 THEN IAK:=HMKHHEPAVVC;
IF (J/(5*(L+I1)+1820))<0276 THEN N:=N;
K:=1;
ADM:DECLARE OP LABEL;
XBJ:IF ((N)+(K+14+((2)+(2+4/04))-(3)+(4)/(6383))+500)+82*(N+83)+2500/(K+1/0)<(1) THEN GOTO JK;
IF 2>32 THEN IF (K+1)<0 THEN 1:=2+4 THEN IF 1>(1K+0) THEN IF K>L THEN IVM:=(K)*1+L;
END SCS;

IJIVIAE:PROCEDURE;
E:DECLARE N LABEL;
GOTO Z5;
GOTO KEVL;
Z5:IF (06+6)>1 THEN OSM=ND;
ND:=I1 Q;
LGFATQ=M;
DECLARE Q LABEL;
RIGI:DECLARE OM LABEL;
TOD:DECLARE R LABEL;
IF 6<11 THEN OWE=8;
TO:DECLARE T LABEL;
IF 6<15 THEN IF 5=ML THEN 4 GOTO OWE;
DECLARE OSM LABEL;
T=R;
X:IF 1>3 THEN GOTO T;
B=X;
MKK=6;
GOTO OSM;
WEVL:GOTO M;
J:DECLARE R LABEL;
END IJIVIAE;

X:PROCEDURE;
IF 915=3 THEN XD=SV;
MQD:GOTO F;
A:=FF=A;
U:DECLARE F LABEL;
Z:IF 98<00274 THEN IF (67+(N-1((KN))+*1+7)-1*(N+39776))<1 THEN GOTO A;
RUMB:DECLARE JFEO LABEL;
DECLARE C LABEL;
SV:=MQ+N+((L+M)*4);
GT:DECLARE K LABEL;
N=MK+6;
C=K;
PRLD:DECLARE FF LABEL;
E:IF 0>MCWC THEN L=N;
GIV:DECLARE S LABEL;
J:DECLARE WYKEL LABEL;
UVGT:DECLARE H LABEL;
EV:=FF=SV;
KS:DECLARE XD LABEL;
DE:DECLARE IC LABEL;
END X;

```

in place of formulations with direct recursion (e.g., the original statements 2220, 2230).

Table 3 gives part of the output obtained from running the system with Little PL/I updated by the above statements.

:recursion limit: <arithmetic expression> 2 3000

In the 10 produced programs, the nesting of arithmetic expressions within arithmetic expressions will be limited to a depth of 2.

:rewrite trace: 3001

The system will "trace rewriting" by printing each nonterminal that is rewritten and the rule selected for it.

:print generated rules: 3001

The system will "trace rule generation" by printing all rules that are dynamically created.

:backtrack trace: 3001

The system will "trace backtracking" by printing all nonterminals whose rewriting directly causes backtracking, with the reason for this backtracking.

:expansion depth: 2 3001

:expansion dictionary: each <boolean expression> 3002

:rewrite systematic: <+ pl/i program> 10 4000

The system will be put into the "systematic mode" and will produce a set of at most 10 programs in which every possible form of Boolean expression occurs in every possible context. The expansion depth statement (together with the language definition) defines when two Boolean expressions are to be considered to have different form.

:error group: <arithmetic expression>, <boolean expression> 3001

:error percent: 50 3002

:error dictionary: all <boolean expression> 3003

:rewrite: <+ pl/i program> 10 4000

In the 10 produced programs, some Boolean expressions will be invalidly rewritten as arithmetic expressions. The first statement says that arithmetic expressions and Boolean expressions are to be considered invalid forms of each other. The second statement says that fifty percent of Boolean expressions are to be invalid, and the third statement says that invalid Boolean expressions may occur in all possible contexts. The produced programs could be used as diagnostic test cases.

## Conclusion

The machine production of programs for testing certain aspects of programming products has been achieved. The system has been successfully used on a number of products to establish their reliability in accepting new test cases without error. The input to the system is a syntax definition in a formal notation. The construction of such a definition for a high-level language is an exacting task. It yields important bonuses by deepening knowledge of the structure of the language and showing up obscurities or ambiguities in the existing documentation. Definitions exist for ECMA Algol, FORTRAN IV, and a major subset of PL/I. The future effort required to adapt the definitions to particular versions of these languages should be small.

## ACKNOWLEDGEMENTS

The PL/I implementation of the syntax machine was designed and programmed by S. M. Glassover, K. V. Hanford, and C. B. Jones.

The systematic and error modes were added by P. D. Wright. P. D. Wright was also responsible for the design and implementation of the assembler version, with the assistance of Miss B. Conyers, a vacation student from Brighton College of Technology, and Mrs. J. Moss.

#### FOOTNOTES AND CITED REFERENCES

1. The syntax machine was implemented for use on System/360 computers.
2. J. W. Backus, "The syntax and semantics of the proposed international algebraic language of the Zürich ACM-GAMM Conference," *Proceedings of the International Conference on Information Processing*, 125-132 (June 1959).
3. N. Chomsky, "Three models for the description of language," *I.R.E. Transactions on Information Theory* **IT-2**, *Proceedings of the Symposium on Information Theory* (September 1956).
4. J. J. Donovan and H. F. Ledgard, "A formal system for the specification of the syntax and translation of computer languages," *AFIPS Conference Proceedings, Fall Joint Computer Conference* **31**, 553-580 (1967).
5. The input philosophy and the structure of grammars in memory were influenced by the compiler-compiler of Brooker et al.<sup>6</sup>
6. R. A. Brooker et al., "The compiler-compiler," R. Goodman, editor, *Annual Review in Automatic Programming* **3**, Pergamon Press, New York (1963).