# A Nonrecursive Method of Syntax Specification

JOHN W. CARR III AND JEROME WEILAND
*The Moore School, University of Pennsylvania, Philadelphia,*
  *Pennsylvania*

The use of the Kleene regular expression notation for describing algebraic language syntax, in particular of ALGOL, is described in this paper. A FORTRAN II computer program for carrying out the elimination algorithm of Gorn, similar to Gaussian elimination for linear systems of algebraic equations, is described. This was applied to numerous smaller languages, including some sublanguages of ALGOL. A hand calculation result of the application of the algorithm to all of ALGOL is given, thus expressing the Revised ALGOL 1960 syntax in completely nonrecursive terms, as far as its context-free portion is concerned. This description in many ways is far more intuitively understood than the previous recursive description, it is suggested. The paper also includes results of the machine program, which does not include a simplification algorithm.

The basis and the method to produce a new approach to computer language syntax specification is outlined in this paper. Given a recursive specification for a context-free language in standard so-called "Backus Normal Form" (BNF) [1] a nonrecursive specification can be produced by using the elimination algorithm of Gorn [2]. The elimination algorithm will solve a set of "equations" (i.e., productions) in BNF in a way similar to that of the standard Gaussian elimination. The elimination algorithm will remove all recursion only on linear languages or Chomsky Type 3 languages [3]. (By "linear" we actually mean "one-sided linear.") If the language is not linear the recursion in the linear portions of the language can be removed, thus producing an overall reduction in recursion in the specification of the language. The basis of the elimination algorithm is:

$$\langle a \rangle ::= \langle a \rangle \langle b \rangle \mid \langle c \rangle$$

is replaced by

$$\langle a \rangle ::= \langle c \rangle (\langle b \rangle) *$$

Here the asterisk indicates zero or more occurrences but not necessarily the same strings drawn from the set $\langle b \rangle$.

After applying the elimination algorithm to a linear language the nonrecursive specification produced is in the form of a regular expression as described by Kleene [4]. A program has been written in FORTRAN II which will take as input a BNF specification for a linear language and produce as output the regular expression [5]. An inverse process has been programmed by Roberts [6]. That context-free portion of ALGOL as defined by Naur [7] was used as input to the program. Upon using Section 2.5.1 ($\langle$number$\rangle$) of Naur as input, a sample of the output from the program is shown in Figure 1.

If the language input to the program is nonlinear the program will go through the elimination procedure an equation (production) at a time until a nonlinear equation is reached. Then the program produces results obtained up to this equation and then stops. The results of the program on ALGOL as well as other hypothetical languages are given by Weiland [5].

The output of the program demonstrates the fact that the regular expressions produced at the end of the back substitution are large and unwieldy to work with except in a computer with a very large memory. The work of Iverson [8] inspired a compact nonrecursive specification of the context-free portion of ALGOL which is presented in Table I. Table I was produced by a hand calculation rather than the machine program, since the program does not include efficient simplification algorithms. Back substitution has not been carried out in general in Table I, in order to keep the specifications compact. Since each equation (production) as originally formulated is im-

```
U I = D I ( D I ) * ,

I N = ( L / 0 / 1 ) D I ( D I ) * ,

D F = 2 D I ( D I ) * ,

E P = 1 0 ( L / 0 / 1 ) D I ( D I ) * ,

D N = ( L / D I ( D I ) * ) 2 D I ( D I ) * / L D I ( D
I ) * ) ,

U N = ( L / ( ( L / D I ( D I ) * ) 2 D I ( D I ) * / L
D I ( D I ) * ) ) 1 0 ( L / 0 / 1 ) D I ( D I ) * / L (
( L / D I ( D I ) * ) 2 D I ( D I ) * / L D I ( D I ) *
) ) ,

N U = ( L / 0 / 1 ) ( ( L / ( ( L / D I ( D I ) * ) 2 D
I ( D I ) * / L D I ( D I ) * ) ) 1 0 ( L / 0 / 1 ) D I
( D I ) * / L ( ( L / D I ( D I ) * ) 2 D I ( D I ) * /
L D I ( D I ) * ) ) ,
```

FIG. 1

## TABLE I

| Name | Nonrecursive Specification |
|---|---|
| **1.1** EMPTY STRING | $\varphi ::=$ |
| **2.1** LETTER | $\langle LE \rangle ::= a|b|c| \cdots |z|A|B|C| \cdots |Z$ |
| **2.2.1** DIGIT | $\langle DI \rangle ::= 0|1|2|3|4|5|6|7|8|9$ |
| **2.2.2** LOGICAL VALUE | $\langle LV \rangle ::=$ **true**\|**false** |
| **2.3** DELIMITERS | (identical with ALGOL 1960 Revised Report) |
| **2.4.1** IDENTIFIER | $\langle ID \rangle ::= LE \ (LE|DI)*$ |
| **2.5.1** UNSIGNED INTEGER | $\langle UI \rangle ::= \langle DI \rangle \langle DI \rangle *$ |
| INTEGER | $\langle IN \rangle ::= (\varphi|+|-)\langle UI \rangle$ |
| DECIMAL FRACTION | $\langle DF \rangle ::= .\langle UI \rangle$ |
| EXPONENT PART | $\langle EP \rangle ::= {}_{10}\langle IN \rangle$ |
| DECIMAL NUMBER | $\langle DN \rangle ::= (\varphi|\langle UI \rangle)\langle DF \rangle|\langle UI \rangle$ |
| UNSIGNED NUMBER | $\langle UN \rangle ::= \langle DN \rangle|(\varphi|\langle DN \rangle)\langle EP \rangle$ |
| NUMBER | $\langle NU \rangle ::= (\varphi|+|-|)\langle UN \rangle$ |
| **2.6.1** PROPER STRING | $\langle PS \rangle ::= \langle AN \rangle|\varphi$ |
| ANY SEQUENCE OF BASIC SYMBOLS NOT CONTAINING "OR" | $\langle AN \rangle$ |
| OPEN STRING | $\langle OS \rangle ::= \langle PS \rangle|'\langle OS \rangle'|\langle OS \rangle \langle OS \rangle$ |
| STRING | $\langle ST \rangle ::= '\langle OS \rangle'$ |
| **3.** EXPRESSION | $\langle EX \rangle ::= \langle BE \rangle|\langle AE \rangle|\langle DE \rangle$ |
| **3.1.1** VARIABLE IDENTIFIER | $\langle VI \rangle ::= \langle ID \rangle$ |
| SIMPLE VARIABLE | $\langle SV \rangle ::= \langle VI \rangle$ |
| SUBSCRIPT EXPRESSION | $\langle SE \rangle ::= \langle AE \rangle$ |
| SUBSCRIPT LIST | $\langle SL \rangle ::= \langle SE \rangle(,\langle SE \rangle)*$ |
| ARRAY IDENTIFIER | $\langle AI \rangle ::= \langle ID \rangle$ |
| SUBSCRIPTED VARIABLE | $\langle SR \rangle ::= \langle AI \rangle[\langle SL \rangle]$ |
| VARIABLE | $\langle VA \rangle ::= \langle SV \rangle|\langle SR \rangle$ |
| **3.2.1** PROCEDURE IDENTIFIER | $\langle PI \rangle ::= \langle ID \rangle$ |
| ACTUAL PARAMETER | $\langle AP \rangle ::= \langle ST \rangle|\langle EX \rangle|\langle AI \rangle|\langle SW \rangle|\langle PI \rangle$ |
| LETTER STRING | $\langle LS \rangle ::= \langle LE \rangle \langle LE \rangle *$ |
| PARAMETER DELIMITER | $\langle PD \rangle ::= ,|)\langle LS \rangle:($ |
| ACTUAL PARAMETER LIST | $\langle AL \rangle ::= \langle AP \rangle(\langle PD \rangle \langle AP \rangle)*$ |
| ACTUAL PARAMETER PART | $\langle AT \rangle ::= \varphi|(\langle AL \rangle)$ |
| FUNCTION DESIGNATOR | $\langle FD \rangle ::= \langle PI \rangle \langle AT \rangle$ |
| **3.3.1** ADDING OPERATOR | $\langle AO \rangle ::= +|-$ |
| MULTIPLYING OPERATOR | $\langle MO \rangle ::= \times|/|\div$ |
| PRIMARY | $\langle PR \rangle ::= \langle UN \rangle|\langle VA \rangle|\langle FD \rangle|(\langle AE \rangle)$ |
| FACTOR | $\langle FT \rangle ::= \langle PR \rangle(\uparrow\langle PR \rangle)*$ |
| TERM | $\langle TE \rangle ::= \langle FT \rangle((\langle MO \rangle \langle FT \rangle))*$ |
| SIMPLE ARITHMETIC EXPRESSION | $\langle SA \rangle ::= ((\varphi|\langle AO \rangle)\langle TE \rangle)((\langle AO \rangle \langle TE \rangle))*$ |
| IF CLAUSE | $\langle IC \rangle ::=$ **if** $\langle BE \rangle$ **then** |
| ARITHMETIC EXPRESSION | $\langle AE \rangle ::= (\langle IC \rangle \langle SA \rangle$ **else**$)*\langle SA \rangle$ |
| **3.4.1** RELATIONAL OPERATOR | $\langle RO \rangle ::= <|\leq|=|\geq|>|\neq$ |
| RELATION | $\langle RE \rangle ::= \langle SA \rangle \langle RO \rangle \langle SA \rangle$ |
| BOOLEAN PRIMARY | $\langle BP \rangle ::= \langle LV \rangle|\langle VA \rangle|\langle FD \rangle|\langle RE \rangle|(\langle BE \rangle)$ |

| Name | Nonrecursive Specification |
|---|---|
| BOOLEAN SECONDARY | $\langle BS \rangle ::= (\varphi|\neg)\langle BP \rangle$ |
| BOOLEAN FACTOR | $\langle BF \rangle ::= \langle BS \rangle(\wedge\langle BS \rangle)*$ |
| BOOLEAN TERM | $\langle BT \rangle ::= \langle BF \rangle(\vee\langle BF \rangle)*$ |
| IMPLICATION | $\langle IP \rangle ::= \langle BT \rangle(\supset\langle BT \rangle)*$ |
| SIMPLE BOOLEAN | $\langle SB \rangle ::= \langle IP \rangle(\equiv\langle IP \rangle)*$ |
| BOOLEAN EXPRESSION | $\langle BE \rangle ::= (\langle IC \rangle \langle SB \rangle$ **else**$)*\langle SB \rangle$ |
| **3.5.1** LABEL | $\langle LA \rangle ::= \langle ID \rangle|\langle UI \rangle$ |
| SWITCH IDENTIFIER | $\langle SW \rangle ::= \langle ID \rangle$ |
| SWITCH DESIGNATOR | $\langle SG \rangle ::= \langle SW \rangle[\langle SE \rangle]$ |
| SIMPLE DESIGNATIONAL EXPRESSION | $\langle SX \rangle ::= \langle LA \rangle|\langle SG \rangle|(\langle DE \rangle)$ |
| DESIGNATIONAL EXPRESSION | $\langle DE \rangle ::= (\langle IC \rangle \langle SX \rangle$ **else**$)*\langle SX \rangle$ |
| **4.1.1** UNLABELLED BASIC STATEMENT | $\langle UB \rangle ::= \langle AS \rangle|\langle GS \rangle|\langle DS \rangle|\langle PT \rangle$ |
| BASIC STATEMENT | $\langle BA \rangle ::= ((\langle LA \rangle:)*\langle UB \rangle$ |
| UNCONDITIONAL STATEMENT | $\langle US \rangle ::= \langle BA \rangle|\langle CS \rangle|\langle BL \rangle$ |
| STATEMENT | |
| COMPOUND TAIL | $\langle SM \rangle ::= \langle US \rangle|\langle CD \rangle|\langle FS \rangle$ |
| BLOCK HEAD | $\langle CT \rangle ::= ((\langle SM \rangle;)*(\langle SM \rangle$ **end**$)$ |
| UNLABELLED COMPOUND | $\langle BH \rangle ::=$ **begin** $\langle DC \rangle(;\langle DC \rangle)*$ |
| UNLABELLED BLOCK | $\langle UC \rangle ::=$ **begin** $\langle CT \rangle$ |
| COMPOUND STATEMENT | $\langle UL \rangle ::= \langle BH \rangle;\langle CT \rangle$ |
| BLOCK | $\langle CS \rangle ::= ((\langle LA \rangle:)*\langle UC \rangle$ |
| PROGRAM | $\langle BL \rangle ::= ((\langle LA \rangle:)*\langle UL \rangle$ |
| | $\langle PR \rangle ::= \langle BL \rangle|\langle CS \rangle$ |
| **4.2.1** LEFT PART | $\langle LP \rangle ::= \langle VA \rangle := |\langle PI \rangle :=$ |
| LEFT PART LIST | $\langle LL \rangle ::= \langle LP \rangle \langle LP \rangle *$ |
| ASSIGNMENT STATEMENT | $\langle AS \rangle ::= \langle LL \rangle(\langle AE \rangle|\langle BE \rangle)$ |
| **4.3.1** GO TO STATEMENT | $\langle GS \rangle ::=$ **go to** $\langle DE \rangle$ |
| **4.4.1** DUMMY STATEMENT | $\langle DS \rangle ::= \varphi$ |
| **4.5.1**[1] IF STATEMENT | $\langle IS \rangle ::= \langle IC \rangle \langle US \rangle$ |
| CONDITIONAL STATEMENT | $\langle CD \rangle ::= ((\langle LA \rangle:)*(\langle IS \rangle(\varphi|$**else** $\langle SM \rangle|\langle IC \rangle \langle FS \rangle))$ |
| **4.6.1** FOR LIST ELEMENT | $\langle FE \rangle ::= \langle AE \rangle(\varphi|$**while** $\langle BE \rangle|$ **step** $\langle AE \rangle$ **until** $\langle AE \rangle)$ |
| FOR LIST | $\langle FR \rangle ::= \langle FE \rangle(,\langle FE \rangle)*$ |
| FOR CLAUSE | $\langle FC \rangle ::=$ **for** $\langle VA \rangle ::= \langle FR \rangle$ **do** |
| FOR STATEMENT | $\langle FS \rangle ::= ((\langle LA \rangle:)*(\langle FC \rangle \langle SM \rangle)$ |
| **4.7.1** PROCEDURE STATEMENT | $\langle PT \rangle ::= \langle PI \rangle \langle AT \rangle$ |
| **5.** DECLARATION | $\langle DC \rangle ::= \langle TD \rangle|\langle AD \rangle|\langle SD \rangle|\langle PA \rangle$ |
| **5.1.1** TYPE LIST | $\langle TP \rangle ::= (\langle SV \rangle,)*\langle SV \rangle$ |
| TYPE | $\langle TY \rangle ::=$ **real**\|**integer**\|**Boolean** |
| LOCAL OR OWN TYPE | $\langle OW \rangle ::= (\varphi|$**own**$) \ \langle TY \rangle$ |
| TYPE DECLARATION | $\langle TD \rangle ::= \langle OW \rangle \langle TP \rangle$ |
| **5.2.1** LOWER BOUND | $\langle LB \rangle ::= \langle AE \rangle$ |
| UPPER BOUND | $\langle UR \rangle ::= \langle AE \rangle$ |
| BOUND PAIR | $\langle BD \rangle ::= \langle LB \rangle:\langle UR \rangle$ |
| BOUND PAIR LIST | $\langle BR \rangle ::= \langle BD \rangle(,\langle BD \rangle)*$ |
| ARRAY SEGMENT | $\langle AY \rangle ::= ((\langle AI \rangle,)*((\langle AI \rangle[\langle BR \rangle])$ |
| ARRAY LIST | $\langle AA \rangle ::= \langle AY \rangle(,\langle AY \rangle)*$ |

[1] IF clause and unconditional statement are not repeated here as in report of NAUR [6].

TABLE I (Continued)

| Name | Nonrecursive Specification |
|---|---|
| ARRAY DECLARATION | $\langle AD \rangle ::= (\textbf{array}|\langle\langle OW \rangle \textbf{ array }) \\ \langle AA \rangle$ |
| 5.3.1 |  |
| SWITCH LIST | $\langle SH \rangle ::= \langle DE \rangle(,\langle DE \rangle)*$ |
| SWITCH DECLARATION | $\langle SD \rangle ::= \textbf{switch } \langle SW \rangle := \langle SH \rangle$ |
| 5.4.1 |  |
| FORMAL PARAMETER | $\langle FM \rangle ::= \langle ID \rangle$ |
| FORMAL PARAMETER LIST | $\langle FA \rangle ::= \langle FM \rangle(\langle PD \rangle\langle FM \rangle)*$ |
| FORMAL PARAMETER PART | $\langle FO \rangle ::= \varphi|(\langle FA \rangle)$ |
| IDENTIFIER LIST | $\langle IF \rangle ::= \langle ID \rangle(,\langle ID \rangle)*$ |
| VALUE PART | $\langle VP \rangle ::= \textbf{value } \langle IF \rangle;|\varphi$ |
| SPECIFIER | $\langle SP \rangle ::= \textbf{string}|\textbf{array}|\langle TY \rangle \\ (\varphi|\textbf{array}|\textbf{proce-} \\ \textbf{dure})|\textbf{label}| \\ \textbf{switch}|\textbf{procedure}$ |
| SPECIFICATION PART | $\langle SF \rangle ::= (\varphi|\langle SP \rangle\langle IF \rangle:)(\langle SP \rangle \\ \langle IF \rangle:)*$ |
| PROCEDURE HEADING | $\langle PH \rangle ::= \langle PI \rangle\langle FP \rangle; \langle VP \rangle\langle SF \rangle$ |
| PROCEDURE BODY | $\langle PO \rangle ::= \langle SM \rangle|\langle CO \rangle$ |
| PROCEDURE DECLARATION | $\langle PA \rangle ::= (\varphi|\langle TY \rangle)(\textbf{procedure} \\ \langle PH \rangle\langle PO \rangle)$ |

mediately solvable in terms of the asterisk notation, a nonrecursive regular expression is produced. This is a special characteristic, not emphasized up until now, of the published context-free part of ALGOL (as is well known for context-free languages in general).

The section numbers in Table I refer to section numbers in the ALGOL report of Naur [7]. Each string class name of the original report is represented in this Table by a symbol composed of two letters. The reason for choosing two letters was that the program was written to accept two characters for each name and that letters could be used as mnemonics for the actual names. No duplicate symbols were allowed; thus the symbol for the name is not always clearly mnemonic. By the use of the distributive law of concatenation over set union the nonrecursive equations have been factored if possible. The names that are defined recursively in the original ALGOL report [7] are easily recognized by the appearance of a * in the nonrecursive specification.

Since the nonrecursive specification requires parentheses as control characters, whenever right or left parenthesis denote themselves in Table I they are in boldface. Section 4.7.1 was not included in its entirety as in [7] because ⟨Procedure statement⟩ has the same specification as 3.2.1 ⟨Function Designator⟩. The symbol $\varphi$ is used instead of ⟨empty⟩ to conform more closely to the Kleene notation. In addition, Section 2.3 (Delimiters) has been omitted as an exact duplicate of the report.

In Figure 1, because of the limitations of a standard computer printer, the symbol $L$ has been used instead of $\varphi$, 0 for +, 1 for −, and 2 for . .

RECEIVED NOVEMBER, 1965

---

## George E. Forsythe, Editor of a
### New Education Department
### Invites Contributions

Computing and Education come together in two different ways. First, the digital computer can be programmed into a powerful tool in the educational process itself, for example as a teaching machine or as a processor of records about student progress. We might call this *computers in education*, and we welcome contributions in this area. (If they deal with educational data-processing techniques that are mainly the same as data-processing techniques for other purposes, articles should be directed to another department of *Communications*.) The second confluence of Computing and Education might be called *education in computing*. This deals with matters of curriculum, personnel, and organization in formal education at all levels about the computing and information sciences. We welcome contributions in this area also. Reference [1] is an excellent preliminary report on education in computing, and it is criticized in reference [3]. Reference [2] deals with both our subjects, discussing a use of *computers in education about computing*.

With the great growth of interest in teaching machines and the sudden emergence of numerous university departments of computer science (under various titles), there should be a great deal of valuable material for this department. Let's have it!—G. E. FORSYTHE

*REFERENCES*

1. ACM Curriculum Committee on Computer Science. An undergraduate program in computer science—preliminary recommendations, *Comm. ACM 8* (Sept. 1965), 543–552.
2. FORSYTHE, GEORGE E. AND WIRTH, NIKLAUS. Automatic grading programs. *Comm. ACM 8* (May 1965), 275–278.
3. PARNAS, DAVID L. On the preliminary report of C³S (letter to the Editor) *Comm. ACM 9* (Apr. 1966), 242–243.

◆

## REFERENCES

1. BACKUS, J. W. The syntax and semantics of the proposed international algebraic language of the Zurick ACM-GAMM conference: ICIP. Paris, June 1959.
2. GORN, S. Processors for infinite codes of the Shannon-Fano type. Proc. of a Symp. on Mathematical Theory of Automata, Polytechnic Institute of Brooklyn, 1962.
3. CHOMSKY, N. On certain formal properties of grammars. *Inf. Contr. 2* (June 1959), 137–167.
4. KLEENE, S. C. *Representations of Events in Nerve Nets and Finite Automata*. Automata Studies, Shannon, C., and McCarthy, J. (Eds.). Princeton, 1956.
5. WEILAND, J. N. Applications of the elimination algorithm to recursive language specifications. Master's Thesis, U. of Pennsylvania, May 1965.
6. ROBERTS, M. B. A generalized recognizer for finite state languages. Master's Thesis, U. of Pennsylvania, August 1965.
7. NAUR, P. (Ed.) Revised report on the algorithmic language ALGOL 60. *Comm. ACM 6*, 1 (Jan. 1963), 1–17.
8. IVERSON, K. E. A method of syntax specification. *Comm. ACM 7*, 10 (Oct. 1964), 588–589.