

PUC

Series: Monographs in Computer Science
and Computer Applications

Nº 10/72

ON THE REDUCTION OF THE SIZE OF TRANSITION MATRICES

by

S.R.P. Teixeira

and

L.F.A. Cunha

Computer Science Department - Rio Datacenter

Pontifícia Universidade Católica do Rio de Janeiro
Rua Marquês de São Vicente, 209 — ZC-20
Rio de Janeiro — Brasil

ON THE REDUCTION OF THE SIZE OF TRANSITION MATRICES

S.R.P. Teixeira
Associate Professor

and

L.F.A. Cunha
Assistent Professor

Computer Science Department
PUC/RJ

Series Editor: Prof. A. L. Furtado

November/1972

ABSTRACT

Although the general context free grammar has shown to be a suitable model for programming language, it has not been used in practice due to time and memory space limitations. More restricted models which guarantee parsing time proportional to n (where n is the length of the input string) have been used instead. Among these schemes, the use of transition matrices has proved to be remarkably efficient. The great speed of this method lies in the fact that for each table reference it does, it makes a syntactic reduction on the sentential form, no further searching being necessary.

The transition matrix is accessed at each step of compilation given the element at the top of the stack and the next element on the input string. Some entries in the matrix specify reductions to be made while others specify error conditions.

The main disadvantage of the method is the large size of the transition matrices for practical programming languages.

This paper presents a method to partition the original grammar into several grammars so that several transition matrices are used instead of just one. Sufficient conditions are given to allow the processor to switch matrices as it is parsing a sentence of the original grammar. This method has all the advantages of the transition matrices technique (mainly speed), while sensibly reducing the memory space required (the most serious drawback of the original method).

Besides, for some grammars and partitions the new method is more powerful than the original method, as shown.

The method presented in this paper is being used in the parser of a fast resident PL/I compiler for the IBM 370/165, being developed at Pontificia Universidade Católica do Rio de Janeiro. It has shown great speed and moderate memory requirements.

1. INTRODUCTION

The use of formal syntax allows for the algorithmic construction of mechanical analyzers (or recognizers) for a significant class of languages. Although some available algorithms encompass a large class of these languages, their practical use has been restricted to somewhat particular cases. This restriction has been imposed by storage medium and program execution time limitations in practical compilers. For this reason, although the general context free has shown to be a suitable model for programming language, it has not been used in practice.

The most efficient algorithm known to date for parsing a general context free language [1], guarantees time proportional to at most n^2 , where n is the length of the string being parsed (in the case of unambiguous grammars).

In practice more restricted models which guarantee time proportional to n have been used. Among these schemes, the use of transition matrices [2] has proved to be remarkably efficient. The great speed of this method lies in the fact that for each table reference it does, it makes a syntax reduction on the input string, no further searching being necessary. Besides, it allows for good error detection

too. These properties are consequence of the very nature of a transition matrix. The matrix is accessed at each step of compilation given the element at the top of the stack and the next symbol on the input string. Some entries in the matrix specify reductions to be made. Others specify error conditions.

The class of languages parsable by processors using transition matrices is a subclass of (1,1) bounded context languages^[3]. The outline of an algorithm to test whether a grammar generates a language which belongs to this subclass is given in [2].

The main disadvantage of the method is the large size of the transition matrices for practical programming languages. This paper presents a method to partition the original grammar into several grammars so that several transition matrices are used instead of a single one. This may result in a great reduction of the total storage area used, depending on the partition chosen. The present method also increases the power of the original method.

2. BASIC DEFINITIONS

An alphabet is any finite nonempty set. For the alphabet V , V^* is the set of all strings of finite-length (words) over V , including the null word Λ . V^+ is defined as $V - \{\Lambda\}$.

If X and Y are sets of words, $X.Y = \{\alpha \beta \mid \alpha \in X, \beta \in Y\}$ where $\alpha \beta$ is the concatenation of α and β , $X^0 = \{\Lambda\}$ and $X^{i+1} = X^i . X$, for $i = 0, 1, 2, \dots$

A context free grammar is a 4-tuple $G = (V_N, V_T, S, P)$

where:

- (i) V_N is the alphabet of nonterminal symbols
- (ii) V_T is the alphabet of terminal symbols, $V_N \cap V_T = \phi$
- (iii) $S \in V_N$ is the start symbol of G .
- (iv) P is a set of productions of the form $A \rightarrow \alpha$,
 $A \in V_N$, $\alpha \in V^+$, where $V = V_N \cup V_T$. $A \rightarrow \alpha$
 is an A-production.

Unless otherwise specified, uppercase Latin letters (A,B,C,...) are used for nonterminal symbols, lower case Latin letters at the beginning of the alphabet (a, b, c, ...) for terminal symbols, lower case Latin letters at the end of the alphabet (t, u, v, ...) for

words in V_T^* , lower case Greek letters $(\alpha, \beta, \gamma, \dots)$ for words in V^* .

For $\alpha, \beta \in V^+$, we say $\alpha \xrightarrow{G} \beta$, or simply $\alpha \Longrightarrow \beta$ when G is understood, if $\alpha = \alpha_1 A \alpha_2$, $\beta = \alpha_1 \gamma \alpha_2$, $\alpha_1, \alpha_2 \in V^*$ $A \rightarrow \gamma \in P$, $\alpha \xrightarrow{R} \beta$ if $\alpha \Longrightarrow \beta$ as above and $\alpha_2 \in V_T^*$.

If $\alpha_1 \Longrightarrow \alpha_2 \Longrightarrow \dots \Longrightarrow \alpha_n$, $n > 1$, we say that $\alpha_1 \xrightarrow{+} \alpha_n$, and that $(\alpha_1, \alpha_2, \dots, \alpha_n)$ is a derivation of α_n from α_1 . Similary if $\alpha_1 \xrightarrow{R} \alpha_2 \xrightarrow{R} \dots \xrightarrow{R} \alpha_n$, $n > 1$, then $\alpha_1 \xrightarrow{R+} \alpha_n$ and $(\alpha_1, \alpha_2, \dots, \alpha_n)$ is a canonical derivation of α_n from α_1 . $(\alpha_n, \dots, \alpha_2 \alpha_1)$ is said to be a canonical parse of α_n to α_1 . When α_1 is not mentioned, it is assumed that $\alpha_1 = S$.

Define $\alpha_1 \xrightarrow{*} \alpha_n$ if $\alpha_1 \xrightarrow{+} \alpha_n$ or $\alpha_1 = \alpha_n$. Analogously we define $\alpha_1 \xrightarrow{R*} \alpha_n$.

If $S \xrightarrow{+} \alpha$ ($S \xrightarrow{R+} \alpha$), α is a (canonical) sentential form. If $\alpha \in V_T^*$ then α is a sentence.

G is said to be reduced if for each $A \in V_N$, each A -production $A \rightarrow \gamma \in P$, $S \xrightarrow{*} \alpha_1 A \alpha_2 \Rightarrow \alpha_1 \gamma \alpha_2 \xrightarrow{*} x$ for $\alpha_1, \alpha_2 \in V^*$, $x \in V_T^*$.

The language generated by G , is: $L(G) = \{x / S \xrightarrow{+} x, x \in V_T^*\}$.

For each context free grammar G there is a context free grammar G' (called the reduced form of G) such that $L(G) = L(G')$, and G' is reduced^[4].

If $S \xrightarrow{*} \alpha_1 A \alpha_2$, $A \xrightarrow{+} \beta$, then $S \xrightarrow{+} \alpha_1 \beta \alpha_2 = \gamma$ and β is said to be a phrase (an A-phrase) of γ . In this case β is a maximal A-phrase of γ if it is not properly contained in another A-phrase of γ . β is a simple phrase if $A \Rightarrow \beta$.

G is said to be unambiguous iff each sentence has a unique canonical derivation.

If G is unambiguous, γ is a sentential form of G , and β_1 and β_2 are phrases of γ then β_1 and β_2 have no symbol in common or one includes the other^[4].

Next, we define three functions:

for each $A \in V_N$

$$\text{head}_G(A) = \{a / a \in V_T, \text{ there exists } \alpha \in V^*, A \xrightarrow{+}_G a \alpha\}$$

$$\text{suc}_G(A) = \{X / X \in V, \text{ there exist } \alpha_1, \alpha_2 \in V^*, S \xrightarrow{+}_{G,R} \alpha_1 A X \alpha_2\}$$

$$\text{pred}_G(A) = \{X / X \in V, \text{ there exist } \alpha_1, \alpha_2 \in V^*, S \xrightarrow{+}_{G,R} \alpha_1 X A \alpha_2\}$$

G is an operator grammar iff no production is of the form $A \rightarrow \alpha_1 B C \alpha_2, \alpha_1, \alpha_2 \in V^*, B, C \in V_N$.

3. THE ORIGINAL PARSING METHOD

It is necessary to repeat some of the important results of the method of transition matrices [2] in order to clarify this paper. The results of this section are proved either directly or indirectly in [2].

A context free grammar $G = (V_N, V_T, S_0, P)$ is an augmented operator grammar (A.O.G) iff:

$$(i) \quad V_N = N_* \cup N, \quad N_* \cap N = \phi.$$

N_* is the set of starred nonterminal symbols (SNTS) and N is the set of unstarred nonterminal symbols (UNTS). Elements of N_* will be denoted as A^*, B^*, \dots , elements of N as: A, B, \dots , and elements of $N_* \cup N$ as: $A^{\bar{}}, B^{\bar{}}, \dots$.

A

$$(ii) \quad \perp \in V_T \text{ is an endmarker .}$$

$$(iii) \quad \perp^* \in N_* \text{, is a starred endmarker .}$$

(iv) The only S_0 production in P

$S_0 \rightarrow \perp^* S \perp$, and $\perp^* \rightarrow \perp \in P$ are the only productions in which S_0 , \perp^* or \perp appear.

(v) Each production in P is in one of the forms:

$$A \rightarrow C \quad (3.1)$$

$$A \rightarrow B^* \quad (3.2)$$

$$A \rightarrow B^* C \quad (3.3)$$

$$A^* \rightarrow a \quad (3.4)$$

$$A^* \rightarrow C a \quad (3.5)$$

$$A^* \rightarrow B^* a \quad (3.6)$$

$$A^* \rightarrow B^* C a \quad (3.7)$$

Furthermore for each $A^* \in N_*$ there is at most one A^* -production in P.

The right part α , of an A^* -production $A^* \rightarrow \alpha$ appears as the right part of no other production.

(vi) For each $A \in N$, $A \xrightarrow{+}_G A$ does not hold. If $A \xrightarrow{+}_G A_k$ the sequence $A_0 \xrightarrow{G} A_1 \xrightarrow{G} \dots \xrightarrow{G} A_k$ is unique, $k \geq 1$.

In [5] an algorithm is given to transform any context free grammar into one satisfying (vi) above. In [2] an algorithm is given to transform an operator grammar satisfying (vi) above into an A.O.G., such that the A.O.G. is unambiguous iff the original operator grammar is unambiguous.

Lemma 3.1 - If G is an A.O.G., then:

(a) No sentential form α has the form $\alpha_1 A B \alpha_2$,
 $\alpha_1, \alpha_2 \in V^*$, $A, B \in N$

(b) No sentential form α has the form $\alpha_1 \beta \gamma \alpha_2$,

where β is an A-phrase, γ is a B-phrase,

$\alpha_1, \alpha_2 \in V^*$ $A, B \in N$.

(c) If $\alpha_1 X \gamma \alpha_2$ is a canonical sentential form, and

γ is an A-phrase, $A \in N$ $X \in V$, then $X \in N_*$

(d) If $B \in N$, $x \in V_T^+$ and $B \xrightarrow[R]{+} x$, the last production used in this derivation of B to x is of the form $A^* \rightarrow a$.

If α is a sentential form of an A.O.G., a prime phrase of α is a phrase such that:

(a) It contains at least one terminal symbol or one STNS

(b) It contains no other phrase satisfying (a) other than itself.

When parsing a sentential form of an A.O.G. we will detect and reduce a prime phrase according to the shortest derivation that generates it. This derivation takes one of the forms:

$$A \implies B^* \quad (3.8)$$

$$A \xrightarrow{+} B^* C \text{ for } A \implies B^* D \text{ and } D \xrightarrow{*} C \quad (3.9)$$

$$A^* \implies a \quad (3.10)$$

$$A^* \xrightarrow{+} C a \text{ for } A^* \implies B a \text{ and } B \xrightarrow{*} C \quad (3.11)$$

$$A^* \implies B^* a \quad (3.12)$$

$$A^* \xrightarrow{+} B^* C a \text{ for } A^* \implies B^* D a \text{ and } D \xrightarrow{*} C \quad (3.13)$$

Note that a phrase C , produced by a derivation of the type $A \xrightarrow{+} C$ is not considered as a separate case. This saves steps during the parse since in practice derivations of this kind do not involve any semantic evaluation.

Theorem 3.1 - A canonical sentential form α of A and A.O.G. has one of the forms:

$$(a) \quad B_1^* B_2^* \dots B_\ell^* a_1 a_2 \dots a_m \quad \ell \geq 0, m \geq 1 \quad (3.14)$$

$$(b) \quad B_1^* B_2^* \dots B_\ell^* C a_1 a_2 \dots a_m \quad \ell \geq 1, m \geq 1 \quad (3.15)$$

where $B_1^* = \perp^*$ when $\ell \geq 1$, $a_m = \perp$

Furthemore a leftmost prime phrase of α has one of the forms:

$$B_\ell^*, B_\ell^* C, a_1, C a_1, B_\ell^* a_1, B_\ell^* C a_1$$

When parsing a sentence x of an A.O.G. G we will reach canonical sentential form $\alpha, \alpha \xrightarrow{*}_R x$. The symbols $B_1^*, B_2^*, \dots, B_\ell^*$ are stored in a pushdown stack. B_ℓ^* is the top of the stack. There is a position called E which stores C in case (3.15). For (3.14) $E = \text{empty}$. a_1, a_2, \dots, a_m are the input symbols to be processed yet. G is said to be parsable by a transition matrix processor if from the top of the stack B_ℓ^* , E and a_1 the leftmost prime phrase β of α is uniquely determined, and the nonterminal A^- such that

$$S_0 \xrightarrow{*}_R \alpha_1 A^- \alpha_2 \xrightarrow{+}_R \alpha_1 \beta \alpha_2 = \alpha$$

$$\alpha_1 \in V^*, \quad \alpha_2 \in V_T^*$$

is also uniquely determined. Note that the derivation of A^- to β which makes

$$\alpha_1 A^- \alpha_2 \xrightarrow{+}_R \alpha_1 \beta \alpha_2$$

has one the forms (3.8), (3.9), ..., (3.13). To complete this parsing step, each of the following three symbols: B_ℓ^* , value of E , a_1 , which is part of β is deleted. If $A^- \in N_*$ then it becomes the new top of the stack. Otherwise A^- will be the new value of E . The process is repeated until $E = S_0$.

Note that a sentence of an A.O.G, G , is of the form $\perp \beta \perp$ where β is an S-phrase. The parser starts with \perp^* in the stack, and $\beta \perp$ as the input string to be processed.

A transition matrix M is a matrix which has a row assigned to each SNTS B^* , and a column assigned to each terminal a . So, with the top of the stack B_ℓ^* and the next input symbol a_1 , we will access the element $M_{B_\ell^*, A_1}$ which specifies for each value of E the unique leftmost prime phrase β , and the unique nonterminal A which generates β .

An A.O.G. G is a T-grammar iff Conditions 1 and Condition 2 below hold.

Condition 1 - For each $B^* \in N_*$, $a \in V_T$ at most one of the following three statements is true:

$$(a) \text{ There exists } A \rightarrow B^* \in P, a \in \text{suc}_G(A) \quad (3.17)$$

$$(b) \text{ There exists } A^* \rightarrow B^* a \in P, \quad (3.18)$$

$$(c) \text{ There exists } A^* \rightarrow a \in P, A^* \in \text{suc}_G(B^*) \quad (3.19)$$

Furthermore if (a) holds the UNTS A is unique.

Condition 2 - For each $C \in N$, $B^* \in N_*$, $a \in V_T$ at most one of the following three statements is true:

$$(a) \text{ There exists } A \rightarrow B^*D \in P, a \in \text{suc}_G(D) \quad (3.20)$$

$$D \xrightarrow[G]{*} C$$

$$(b) \text{ There exists } A^* \rightarrow B^*D \in P, D \xrightarrow[G]{*} C \quad (3.21)$$

$$(c) \text{ There exists } A^* \rightarrow D \in P, A^* \in \text{suc}_G(B^*) \quad (3.22)$$

$$D \xrightarrow[G]{*} C$$

Furthermore if (a) holds both A and D are unique. If (b) or (c) hold D is unique.

Before the next Theorem, it is worth to note that the procedure that transforms any context free grammar G into a reduced context free grammar G', $L(G') = L(G)$, when applied to an A.O.G. G, produces another A.O.G. G'.

Theorem 3.2 - If an A.O.G. G is a T-grammar then it is unambiguous.

Theorem 3.3 - If an A.O.G. G is parsable by a transition matrix processor then its reduced form G' is a T-grammar.

Theorem 3.4 - If an A.O.G. G is a T-grammar then it is parsable by a transition matrix processor.

The method of transition matrices is very fast, its disadvantage is the large size of the transition matrix. In the next section we will see how to reduce the size of the matrix.

4. THE PARTITIONING METHOD

Let $G = (V_N, V_T, S_0, P)$ be a reduced A.O.G for which we want to construct a transition matrix processor. We will partition G into two grammars and construct transition matrices for both. When parsing a sentence we will detect conditions which will tell us when to switch matrices. This scheme allows for a sensible reduction of the total memory space occupied by the transition matrices without reducing the speed of the method noticeably.

Let $S_0 \rightarrow \perp^* S \perp$ be the only S_0 -production in P . Let $S_1 \in V_N$, $S \neq S_0$, $S_1 \neq S$.

$H' = (V'_N, V'_T, S_0, \bar{P}')$ is the reduced form of $(V_N, V_T \cup \{t_1\}, S_0, \bar{P})$
 where: $t_1 \notin V$ $\bar{P} = \{A \rightarrow \bar{\alpha} / A \rightarrow \alpha P, A \neq S_1, \bar{\alpha}$ is obtained from α by replacing each occurrence of S_1 by $t_1\}$

$G_0 = (V_{N_0}, V_{T_0}, S_0, P_0)$ is the A.O.G. $(V'_N \cup \{S_1\}, V'_T - \{t_1\}, S_0, P')$

where:

$P' = \{A \rightarrow \underline{\alpha} / A \alpha \in \bar{P}' , \underline{\alpha}$ is obtained from α by replacing each occurrence of t_1 by $S_1\}$.

Note that G_0 is an A.O.G. but it is not reduced.

$H'_1 = (V'_N, V'_T, S_1, P'_1)$ is the reduced form of

$H_1 = (V_N, V_T, S_1, P)$. $G_1 = (V_{N_1}, V_{T_1}, X_1, P_1)$ where:

$X_1 \notin V_N$, $V_{N_1} = V_N \cup \{\underline{1}^*, X_1\}$, $V_{T_1} = V_T \cup \{\underline{1}\}$ $P_1 = P' \cup \{X_1 \rightarrow \underline{1}^* S_1 \underline{1} , \underline{1}^* \rightarrow \underline{1}\}$

There is a complete example in the appendix.

Note that since G_0 and G_1 are A.O.G.'s we will try to construct a transition matrix for each. We will succeed only if both G_0 and G_1 are T-grammars. In case this does not happen the next two lemmas say that G is not a T-grammar.

Lemma 4.1 - If G is a T-grammar then G_0 is a T-grammar.

Proof - The result follows by observing that $P_0 \subseteq P$.

Lemma 4.2 - If G is a T-grammar then G_1 is a T-grammar.

Proof - The proof of this and subsequent Lemmas and Theorems is omitted for lack of space.

The converse of the two previous Lemmas is not true. Both G_0 and G_1 may be T-grammars while G is not. This case is illustrated by the example in the appendix.

We will now proceed to explain the method to detect the conditions to switch tables. Consider the A.O.G. G and assume both G_0 and G_1 are T-grammars. Let $x = a_0 a_1 \dots a_n a_{n+1}$ be a sentence generated by G . We have $a_0 = a_{n+1} = \perp$. Let $\beta = a_{i_1} a_{i_1+1} \dots a_{i_1+n_1-1}$, $n_1 \geq 1$, $i_1 \geq 1$, $i_1 + n_1 - 1 \leq n$ be the leftmost maximal S_1 -phrase in x . Since x is a sentence, β is a prime phrase also. It is essential to consider a maximal phrase because we may have a S_1 -phrase within another S_1 -phrase.

Assume we start a canonical parse of x using a transition matrix processor for G_0 . We will be able to proceed correctly until we reach the canonical sentential form $A_1^* \dots A_j^* a_{i_1} \dots a_{n+1}$ (see Lemma 3.1 (c) and Theorem 3.1), $A_1^* = \perp^*$, $j \geq 1$. Then we will look at the transition matrix for G_0 (denoted as M_0), for the action is the case $E = \text{empty}$, for the pair (A_j^*, a_{i_1}) . Because a S_1 -phrase is about to begin, the correct action in this case should be to reduce the prime

phrase a_{i_1} (see Lemma 3.1 (d)) to the unique A_{j+1}^* , such that $A_{j+1}^* \rightarrow a_{i_1} \in P$. One possibility is that column a_{i_1} may not exist in M_0 (row A_j^* certainly exists in M_0), or in case it does there is no action for (A_j^*, a_{i_1}) and $E = \text{empty}$. In this case we add to M_0 for (A_j^*, a_{i_1}) and $E = \text{empty}$, the action to reduce a_{i_1} to A_{j+1}^* followed by an order to switch (denoted SW) to the transition matrix of G_1 . The other possibility is that there is already an action in M_0 for (A_j^*, a_{i_1}) and $E = \text{empty}$. In this case there is a conflict and we cannot resolve the ambiguity with the context used with the transition matrix. Consequently our method will not work (see Lemma 4.3).

What was said above for a particular pair (A_j^*, a_{i_1}) has to be done for each pair (A^*, a) , where $A^* \in \text{pred}_{G_0}(S_1)$ and $a \in \text{head}_{G_1}(S_1)$. If no conflicts occur we will then obtain a new matrix we will call M_0' .

Now, we want to modify the transition matrix for G_1 (denoted M_1) to allow us to continue the canonical parse correctly until we reach a canonical sentential form as:

$$(a) \quad A_1^* \dots A_j^* S_1 a_{i_1 + n_1} \dots a_{n+1} \quad (4.1)$$

$$(b) \quad A_1^* \dots A_j^* B a_{i_1 + n_1} \dots a_{n+1} \text{ where} \quad (4.2)$$

$$S_1 \xrightarrow[G_1]{+} B$$

We have to consider the case (b) above because productions of the form $A \rightarrow C$ (3.1) do not correspond to prime phrases and consequently are not parsed separately. Note that M_1 will not work above because it uses only \perp as a right delimiter to the maximal S_1 -phrase, and will not use $a_{i_1 + n_1}$.

We say that two elements ξ, ϕ of transition matrices are compatible iff there is no value v of E (including empty) for which the sequence of actions (or orders) in ξ (if any) differs from the sequence of actions in ϕ (if any). In this case, to merge ϕ into ξ will add to ξ all sequences of actions in ϕ which did not exist in ξ . If ξ and ϕ are not compatible, the merging will produce a conflict.

We will now proceed to modify M_1 . For each $A^* \in V_{N_1}$, $A^* \neq \perp^*$ the entry in M_1 for (A^*, \perp) is deleted from M_1 (has all actions erased in M_1) and the original entry is merged into the entry for (A^*, a) for each $a \in \text{succ}_{G_0}(S_1)$ (note that we may have $a = \perp$). Similarly, for each $a \in V_{T_1}$, $a \neq \perp$ the entry in M_1 for (\perp^*, a) is deleted from M_1 , and merged into the entry for (A^*, a) for each $A^* \in \text{pred}_{G_0}(S_1)$. (note that we may have $A^* = \perp^*$). Finally, the entry in M_1 for (\perp^*, \perp) is deleted, and the original entry with all actions substituted for SW's is merged into the entry for (A^*, a) , for

each $A^* \in \text{pred}_{G_0}(S_1)$ $a \in \text{suc}_{G_0}(S_1)$. If no conflicts occur in the above process then the matrix M_1 with all the modifications in called M'_1 .

With M'_1 (instead of M_1) we will reach the situation (4.1) or (4.2) and then we will execute SW and switch back to M'_0 correctly. The return to M'_0 is done at the right time since no two maximal S_1 -phrases may be adjacent (see Lemma 3.1 (b)). If there is any conflict in the construction of M'_1 , then our method will not work (see Lemmas 4.4, 4.5, 4.6).

The transition matrix processor starts with M'_0 and switches between M'_0 and M'_1 whenever it executes an order SW. If no action (or order) is found in one of the matrices, then an error has occurred. This modified processor is called a segmented transition matrix processor

Lemma 4.3 - If there exist $A^* \in \text{pred}_G(S_1)$,
 $a \in \text{head}_G(S_1)$, and in M_0 for (A^*, a)
and $E = \text{empty}$, there is an action which is
different from the action to reduce a to
the unique B^* such that $B^* \rightarrow a \in P$, then
 G is not a T-grammar.

Lemma 4.4 - If there exist $A^* \in V_{N_1}$, $A^* \neq \perp^*$,
 $a \in \text{succ}_{G_0}(S_1)$ such that in M_1 the entry
for (A^*, \perp) is not compatible with the
entry for (A^*, a) then G is not a T-gram-
mar.

Lemma 4.5 - If there exist $a \in V_{T_1}$, $a \neq \perp$,
 $A^* \in \text{pred}_{G_0}(S_1)$ such that in T_1 the entry
for (\perp^*, a) is not compatible with the
entry for (A^*, a) , then G is not a T-gram-
mar.

Lemma 4.6 - If there exist $A^* \in \text{pred}_{G_0}(S_1)$, $a \in \text{succ}_{G_0}(S_1)$
 $A^* \neq \perp^*$, $a \neq \perp$ such that in T_1 the
entry for (\perp^*, \perp) with all actions sub-
stituted for SW orders is not compatible
with the entry for (A^*, a) , then G is not
a T-grammar.

Theorem 4.1 - If M'_0 and M'_1 can be constructed with-
out any conflict then G is unambiguous and the language parsable by
the segmented transition matrix processor is precisely $L(G)$.

A algorithm to produce a segmented transition matrix processor for G using G_0 and G_1 will follow the steps detailed in this paper. First it tests whether G_0 and G_1 are T-grammars in order to construct M_0 and M_1 [2]. Then it constructs M'_0 and M'_1 , if any conflict occurs the method fails. In this case Lemmas 4.3, 4.4, 4.5 and 4.6 provide conditions to determine whether the original grammar G was a T-grammar. In the appendix an example is given where G is not a T-grammar, but nevertheless there is a segmented transition matrix processor for G . Furthermore, the product of the number of SNTS and the number of terminal symbols of G is 154. But M'_0 is only $9 \times 7 = 56$ and M'_1 is $7 \times 6 = 42$. So, the savings of memory positions as compared to 154 is over 30%.

5. EXTENSION OF THE PARTITION METHOD

The method described in the previous section may be extended when the grammar G is partitioned in more than two parts. In this case we will have the grammars G_0, G_1, \dots, G_i with matrices M_0, M_1, \dots, M_i respectively. Analogously to what has been done when $i=1$, we construct matrices M'_0, M'_1, \dots, M'_i . The processor starts with matrix M'_0 and may switch to M'_j , $1 \leq j \leq i$. Then it may switch to M'_k , $1 \leq k \leq i$, $k \neq j$ and so forth. In this general case we will need a pushdown stack to be able to remember the matrix (M'_j above) to return to when we finish processing using M'_k . This was not necessary when $i=1$.

The results, Lemmas, etc... for this general case are analogous to the case when $i=1$, and the details will not be given here.

6. CONCLUSIONS

The method presented has all the advantages of the transition matrices technique (mainly speed), while sensibly reducing the memory space requirement (the most serious drawback of the original method [2]).

Besides the partition of the original grammar may increase the power of the original method in some cases, as shown (appendix).

The method presented in this paper is being used in the parser for a fast resident PL/I compiler for the IBM 370/165, being developed at Pontificia Universidade Católica do Rio de Janeiro. It has shown great speed and moderate memory requirements.

REFERENCES

1. Earley, J. - An Efficient Context Free Parsing Algorithm. Comm. ACM 13 (Feb. 1970), 94-102.
2. Gries, D. - The Use of Transition Matrices in Compiling - Comm. ACM 11 (Jan. 1968), 26-34.
3. Floyd, R.W. - Bounded Context Syntactic Analysis. Comm. ACM 7 (Feb. 1964), 62-67.
4. Hopcroft, J. and Ullman, J. - Formal Languages and their Relation to Automata - Addison Wesley, New Yourk, 1969.
5. Ginsburg, S. - An Introduction to Mathematical Theory of Context - Free Languages. McGraw-Hill Book Company, N.Y., 1966.

APPENDIX

Let N be an operator grammar whose productions appear below:

$$\begin{array}{ll}
 S \rightarrow a S_1 a S_2 & S_2 \rightarrow S_2 e \\
 S_1 \rightarrow C & S_2 \rightarrow e \\
 C \rightarrow C e & S_2 \rightarrow g D h \\
 C \rightarrow e & D \rightarrow i \\
 B \rightarrow B c & D \rightarrow D j \\
 B \rightarrow f & \\
 S_1 \rightarrow b B d &
 \end{array}$$

Using the ideas of the algorithm described in [2] we may obtain the reduced A.O.G. G whose productions appear below:

$$\begin{array}{ll}
 S_0 \rightarrow \perp^* S \perp & \perp^* \rightarrow \perp \\
 A_1^* \rightarrow a & A_2^* \rightarrow e \\
 A_3^* \rightarrow b & A_4^* \rightarrow f \\
 A_5^* \rightarrow i & S \rightarrow A_6^* S_2 \\
 A_6^* \rightarrow A_1^* S_1 a & S_1 \rightarrow C \\
 S_1 \rightarrow A_7^* & A_7^* \rightarrow A_3^* B d
 \end{array}$$

C	→	A ₈ *	A ₈ *	→	C e
C	→	A ₂ *	B	→	A ₉ *
A ₉ *	→	B c	B	→	A ₄ *
S ₂	→	A ₁₀ *	A ₁₀ *	→	S ₂ e
S ₂	→	A ₂ *	S ₂	→	A ₁₂ *
A ₁₂ *	→	A ₁₁ * D h	A ₁₁ *	→	g
D	→	A ₅ *	D	→	A ₁₃ *
A ₁₃ *	→	D _j			

G₀ has 14 SNTS and 11 terminals. Note that G is not a T-grammar.

G₀ has the productions:

S ₀	→	⊥* S ⊥	S	→	A ₆ * S ₂
A ₆ *	→	A ₁ * S ₁ a	A ₁ *	→	a
S ₂	→	A ₁₀ *	A ₁₀ *	→	S ₂ e
S ₂	→	A ₂ *	A ₂ *	→	e
S ₂	→	A ₁₂ *	A ₁₂ *	→	A ₁₁ * D h
A ₁₁ *	→	g	D	→	A ₁₃ *
A ₁₃ *	→	D _j	D	→	A ₅ *
A ₅ *	→	i	⊥*	→	⊥

M₀ and M'₀ are 9 x 7 = 56

G_1 has the productions:

$X_0 \rightarrow \perp^* S_1 \perp$ $C \rightarrow A_8^*$ $A_8^* \rightarrow C e$ $S_1 \rightarrow A_7^*$ $A_3^* \rightarrow b$ $B \rightarrow A_9^*$ $A_9^* \rightarrow B c$	$S_1 \rightarrow C$ $C \rightarrow A_2^*$ $A_2^* \rightarrow e$ $A_7^* \rightarrow A_3^* B d$ $\perp^* \rightarrow \perp$ $B \rightarrow A_4^*$ $A_4^* \rightarrow f$
---	---

M_1 is 7×6 . M_1' does not have a row for \perp^* and a column for \perp but has an extra column for \underline{a} , and an extra row for A_1^* . Consequently we still have that M_1' is 7×6 .