SUBWORDS IN DETERMINISTIC TOL LANGUAGES *

By

A. Ehrenfeucht **
G. Rosenberg ***

** Department of Computer Science
University of Colorado, Boulder, Colorado

*** Department of Computer Science
State University of New York at Buffalo

Report # CU-CS-015-73                    March, 1973

# 1. INTRODUCTION.

Developmental systems are formal structures which model the way in which certain biological organisms develop. The study of such systems has recently attracted particular attention as a new branch of the theory of formal languages, see for example [1], [4], [6] and their references. In this paper we continue research into TOL systems which were introduced in [2] and further studies, for example, in [3], [5].

A TOL system has the following components

(i)   A finite set of symbols, $\Sigma$, the <u>alphabet</u>,

(ii)   A finite collection $P$ of <u>tables</u>, each of which tell us by what strings in $\Sigma^*$ a symbol may be replaced. A table may, in general, contain several productions for each symbol. In every step of a derivation, all symbols in a string must be simultaneously replaced according to the production rules of one arbitrarily chosen table.

(iii) A starting string, $\omega$, the <u>axiom</u>. The language generated by a given TOL system consists of $\omega$ and all strings which can be derived from $\omega$ in a finite number of steps.

A TOL system is called <u>deterministic</u> (abbreviated a DTOL system) if each of its tables is such that for each symbol in the alphabet the table contains exactly one production with the symbol on the left. The role of a deterministic restriction is one of the important questions, from both the biological and formal points of view, in the theory of developmental systems.

In this paper we prove a result, which we believe is a fundamental one for the characterization of languages generated

by DTOL systems (called DTOL languages.)  This result says that if  L  is a DTOL language over an alphabet containing at least two letters then the relative number of subwords of a given length  k  occurring in the words of  L  tends to zero as  k  increases.

The formal definition of a TOL system is given for example in [3] the terminology and notation of which we shall follow in this paper.  In addition to this we shall use the following notation:

1)  If  x  is a word then  $|x|$  denotes its length and if  A  is a set then  #A  denotes the cardinality of  A.  The empty word is denoted by  $\Lambda$.

2)  If  $G = <\Sigma, P, \omega>$  is a DTOL system and  $T \varepsilon P$  then  $a \underset{T}{\rightarrow} \alpha$  abbreviates "the production  $a \rightarrow \alpha$  is in  T", and for a word  x  $T(x)$  denotes the word  y  such that  $x \underset{T}{\Rightarrow} y$.

3)  If  L  is a language and  k  a natural number, then  $P_k(L)$  denotes the set of all subwords of length  k  occurring in the words of  L.

## 2. PRELIMINARY LEMMAS.

Let $\Sigma$ be a finite alphabet where $\#\Sigma = n \geq 2$, and $Z = \{w_1, \ldots, w_n\}$ a finite non-empty subset of $\Sigma^*$ such that $|w| > 1$ for at least one $w$ in $Z$. First three lemmas state some properties of $Z^*$.

### Lemma 1.

There exists $k_0$ such that $P_{k_0}(Z^*) < n^{k_0}$.

Proof:

For $a$ in $\Sigma$, let $W(Z,a) = \{w \varepsilon Z: a \text{ occurs in } w\}$. We can distinguish two cases.

Case I. For every $a$ in $\Sigma$ there exists a word $w$ in $W(Z,a)$ such that $w = a^r$ for some $r \geq 0$.

Case II. There exists a letter $a$ in $\Sigma$ such that no word in $W(Z,a)$ is of the form $a^r$ for some $r \geq 0$.

If Case I holds, then (because $\#Z = \#\Sigma$), for every $a$ in $\Sigma$ there exists exactly one $w$ in $Z$ such that $w = a^r$ for some $r \geq 0$. But then one of the words in $Z$ is of the form $b^\ell$ for some $\ell > 1$, $b$ in $\Sigma$, where $b$ does not occur in any other word in $Z$. Consequently if $x \varepsilon \Sigma - \{b\}$ then $xbx \not\in Z^*$.

If Case II holds then we have to consider two subcases.

Subcase II.1. There exists a letter $a$ in $\Sigma$ such that $W(a,Z) = \emptyset$. Then $a \not\in Z^*$.

Subcase II.2. For every $a$ in $\Sigma$, $W(a,Z) \neq \emptyset$. Then, obviously, $a^{2u_a} \not\in Z^*$ for an arbitrary $a$ in $\Sigma$ and $u_a = \max\{|w| : w \varepsilon W(Z,a)\}$.

Lemma 1 follows now from the above case analysis.

Lemma 2.

Let $k \varepsilon N^+$, $k = k_0 s + k_1$ where $k_0$ is a constant as in Lemma 1 and $k_1 < k_0$. Then

$$\frac{P_k(Z^*)}{n^k} \leq \left(\frac{P_{k_0}(Z^*)}{n^{k_0}}\right)^s .$$

Proof:

Let $k, k_0, s, k_1$ be as described above.
Obviously $P_k(Z^*) \leq (P_{k_0}(Z^*))^s \cdot n^{k_1}$.

Hence

$$\frac{P_k(Z^*)}{n^k} \leq \frac{(P_{k_0}(Z^*))^s \cdot n^{k_1}}{n^{k_0 s + k_1}} = \left(\frac{P_{k_0}(Z^*)}{n^{k_0}}\right)^s$$

and Lemma 2 holds.

Lemma 3.
$$\lim_{k \to \infty} \frac{P_k(Z^*)}{n^k} = 0.$$

Proof:

Let $k_0$ be a constant as in Lemma 1, $k \varepsilon N^+$ and $s$ be defined as in Lemma 2. Then from Lemma 2 it follows that $\frac{P_k(Z^*)}{n^k} \leq \left(\frac{P_{k_0}(Z^*)}{n^{k_0}}\right)^s$ and from Lemma 1 it follows that $\frac{P_{k_0}(Z^*)}{n^{k_0}} < 1$.
Hence

$$\lim_{k \to \infty} \frac{P_k(Z^*)}{n^k} \leq \lim_{s \to \infty} \left(\frac{P_{k_0}(Z^*)}{n^{k_0}}\right)^s = 0$$

and Lemma 3 holds.

For our next Lemma we shall need some additional notation.
Let $G = \langle \Sigma, , \omega \rangle$ be a DTOL system. We shall now distinguish some

subsets of $P$ and accordingly some subsets of $L(G)$.

Let $P_g = \{T \varepsilon P: a \xrightarrow{T} \alpha$ for some a in $\Sigma$, $\alpha$ in $\Sigma*$ such that $|\alpha| \geq 2\}$,

$P_c = \{T \varepsilon P:$ if $a \xrightarrow{T} \alpha$ then $|\alpha| = 1\}$, (each element of $P_c$ is called

a <u>coding table</u>),

$\overline{P}_c = \{T \varepsilon P:$ if $a \xrightarrow{T} \alpha$ then $|\alpha| = 1$ and $b \xrightarrow{T} \Lambda$ for some b in $\Sigma\}$.

If $T \varepsilon P$ then we define two languages "associated with T" as

follows:

$L_T = \{x \varepsilon \Sigma*:$ there exists y in $\Sigma*$ such that $\omega \Rightarrow* y \underset{T}{\Rightarrow} x\}$,

$\hat{L}_T = \{x \varepsilon \Sigma*:$ there exists y in $L_T$ and $\overline{T}$ in $P_c$ such that $y \underset{\overline{T}}{\Rightarrow} x\}$.

Now we shall show how an arbitrary DTOL language can be

decomposed using this notation.

<u>Lemma 4</u>.

For every DTOL system $G$ there exists an equivalent DTOL system

H with a set of tables $P$ such that

$$L(G) = F \cup \bigcup_{T \varepsilon P_g} L_T \cup \bigcup_{T \varepsilon P_g} \hat{L}_T \cup \bigcup_{T \varepsilon \overline{P}_c} L_T$$

where F is a finite language.

Proof:

Let $G = \langle \Sigma, R, \omega \rangle$ be a DTOL system.

We leave to the reader the easy proof of the fact that one may

(effectively) construct a finite number of coding tables $P_1, P_2,$

$\ldots, P_t$ (not necessarily in $R$) such that for every x in $\Sigma*$

and every sequence $T_{i_1} T_{i_2} T_{i_r}$ of coding tables from $R \cup \{P_1, \ldots, P_t\}$

there exists $j \varepsilon \{1, \ldots, t\}$ such that

$$T_{i_r} \ldots T_{i_2} T_{i_1}(x) = P_j(x).$$

Hence we have a DTOL system $H = <\Sigma, P, \omega>$ where
$P = (R - R_c) \cup \{P_1, \ldots, P_t\}$, which is equivalent to $G$.
We leave to the reader the easy proof of the fact that if we set
$F = \{\omega\} \cup \{P(\omega) : P \varepsilon P_c\}$ then indeed

$$L(G) = F \cup \bigcup_{T \varepsilon P_g} L_T \cup \bigcup_{T \varepsilon P_g} \hat{L}_T \cup \bigcup_{T \varepsilon \bar{P}_c} L_T$$

and so Lemma 4 holds.

## 3. MAIN RESULT.

Theorem.

Let $\Sigma$ be a finite alphabet such that $\#\Sigma = n \geq 2$.

If $L$ is a DTOL language, $L \subseteq \Sigma^*$, then

$$\lim_{k \to \infty} \frac{P_k(L)}{n^k} = 0.$$

Proof:

Let $\Sigma$ be a finite alphabet where $\#\Sigma = n \geq 2$ and $L$ be a DTOL language ($L \subseteq \Sigma^*$) generated by a DTOL system $G = \langle \Sigma, P, \omega \rangle$. According to Lemma 4 and its proof we may assume that

$$L = L(G) = F \cup \bigcup_{T \varepsilon P_g} L_T \cup \bigcup_{T \varepsilon P_g} \hat{L}_T \cup \bigcup_{T \varepsilon \overline{P}_c} L_T$$

where $F = \{\omega\} \cup \{T(\omega) : T \varepsilon P_c\}$ is a finite language.

1) Obviously $\lim_{k \to \infty} \dfrac{P_k(F)}{n^k} = 0$.

2) If $T \varepsilon P_g$, then (as $L_T \subseteq \{\alpha : a \underset{T}{\to} \alpha\}^*$) from Lemma 3 it follows that $\lim_{k \to \infty} \dfrac{P_k(L_T)}{n^k} = 0$ and consequently

$$\lim_{k \to \infty} \frac{P_k(\bigcup_{T \varepsilon P_g} L_T)}{n^k} = 0.$$

3) Let $T \varepsilon P_g$. If $\#P_c = m$ then $\#\hat{L}_T \leq m \cdot (\#L_T)$ and consequently $\#P_k(\hat{L}_T) \leq m \cdot P_k(L_T)$. Thus from 2) it follows that $\lim_{k \to \infty} \dfrac{P_k(\hat{L}_T)}{n^k} = 0$, hence

$$\lim_{k \to \infty} \frac{P_k(\bigcup_{T \varepsilon P_g} \hat{L}_T)}{n^k}.$$

4) Let $T \varepsilon \overline{P}_c$. There exists a letter in $\Sigma$ which does not occur at the right hand side of any production in $T$. Thus

$P_k(L_T) \leq (n-1)^k$ and $\dfrac{P_k(L_T)}{n^k} \leq \dfrac{(n-1)^k}{n^k}$, hence $\lim\limits_{k \to \infty} \dfrac{P_k(L_T)}{n^k} = 0$

and consequently

$$\lim_{k \to \infty} \frac{P_k\left(\bigcup\limits_{T \varepsilon \overline{P}_c} L_T\right)}{n^k} = 0.$$

Thus from 1) through 4) and from the expression for $L(G)$ it follows that $\lim\limits_{k \to \infty} \dfrac{P_k(L)}{n^k} = 0$ and so the Theorem holds.

Up to the present, characterization results have been conspicuously absent from the theory, necessitating involved combinatorial proofs to show that certain languages are not DTOL. The result in this paper gives us a direct way to show that many languages are not generable by DTOL systems. For example if $\Sigma = \{a,b\}$ and $F$ is a finite language over $\Sigma$ then $\Sigma^* - F$ is not in the class of DTOL languages.

REFERENCES.

[1]   A. Lindenmayer, G. Rozenberg, Developmental systems and
      languages, in <u>Proc. IVth ACM Symp. on Theory of Comp.</u>,
      Denver, Colorado, 1972.

[2]   G. Rozenberg, T0L systems and languages, in <u>Information and
      Control</u>, to appear.

[3]   G. Rozenberg, The equivalence problem for deterministic T0L
      systems is undecidable, in <u>Information Processing Letters</u>,
      vol. 1, no. 5, 1972.

[4]   G. Rozenberg, D0L sequences, in <u>Discrete Mathematics</u>, to
      appear.

[5]   G. Rozenberg, Extension of tabled 0L systems and languages,
      submitted for publication.

[6]   A. Salomaa, <u>Formal Languages</u>, Academic Press, 1973.