# REPEATABLE  WORDS  FOR  SUBSTITUTION

Taishin   Nishida   &   Youichi   Kobuchi

Department of Biophysics, Kyoto University

## Introduction

Among many language defining mechanisms, sequential rewriting systems, or grammars, and parallel rewriting systems, or L systems, are the two major ways to generate the words of a language systematically.  Many comparative studies of the generative powers between the parallel rewriting systems and the grammars have been investigated.  It is already known that the family of languages generated by 0L systems, i.e., interactionless L systems, and that of context free languages are mutually incomparable [1,6]. In fact very simple 0L systems can generate non-context free languages, e.g., let $S=\langle\{a\}, \tau, a\rangle$ where $\tau$ is a homomorphism on a* given by $\tau(a)=a^2$, then the language generated by S is $L(S)=\{a^{2^i} \mid i\geqslant 0\}$.

Several works have been done to answer the questions why some of the parallel rewriting systems can generate such rather complicated languages, or conversely, which parallel rewriting system merely generates a context free or a rational language. Lindenmayer [4] showed a sufficient condition for a 0L system to generate a context free language.  Král [3] showed a similar condition for an iterated substitution which generates a context free

set. Herman and Walker [2] characterize context free languages with OL systems through "adult" concept; i.e., a word w is adult for a OL system if the descendants of w produced by the OL system consist of w alone. Nishida and Kobuchi [5], and Sakarovitch, Nishida, and Kobuchi [7] introduce a recurrent word which is a generalization of an adult word; i.e., a word w is recurrent if it is a descendant of any of its descendant. Clearly an adult word is recurrent. It is shown in [7] that there is the same characterization of context free languages using iterated substitutions and the recurrent concept.

In [7] it is also shown that the set of recurrent words for a rational (resp. context free) substitution is rational (resp. context free). That is, for any parallel rewriting process, the set of recurrent words is not parallel any more. In this paper we will show that the same statement is valid for the wider set of words which contains the set of recurrent words properly. We first define repeatable word, i.e., a word w is repeatable for a substitution if it is a descendant of itself. Needless to say a recurrent word is repeatable. Then we will show that the set of repeatable words for a rational (resp. context free) substitution is rational (resp. context free). We will show our result with a similar method to that of [7]. Hence we omit the details of some technical lemmas the reader should refer to [7].


1. Preliminaries


Let $\Sigma$ be a nonempty set, called <u>alphabet</u>, the elements of which are called <u>letters</u>. The finite sequences of letters, called

<u>words</u>, together with the operation of concatenation form the free momoid generated by $\Sigma$, $\Sigma^*$ ; the empty word, denoted by 1, is the identity element of $\Sigma^*$.

A subsequence of a word s is called a sparse subword of s or, for short, a <u>subword</u> of s. The <u>length</u> of a word s is, by definition, the length of the sequence s and is denoted by $|s|$ ; if V is any subset of $\Sigma$, $|s|_V$ denotes the number of occurrences of letters of V in s. We denote by $\mathscr{P}(\Sigma)$ the power set of $\Sigma$. The structure of monoid of $\Sigma^*$ extends to the power set $\mathscr{P}(\Sigma^*)$ by $XY=\{xy \mid x \in X \; y \in Y\}$ for any subsets X and Y of $\Sigma^*$. We denote by $card(\Sigma)$ the cardinality of the set $\Sigma$.

A multivalued mapping $\tau : \Sigma^* \to \Theta^*$ is a <u>substitution</u> if it is a homomorphism from $\Sigma^*$ into $\mathscr{P}(\Theta^*)$. Thus a substitution is completely defined by the family of sets $\{\tau(a) \mid a \in \Sigma\}$ and we have $\tau(1)=1$. As any relation, a substitution $\tau : \Sigma^* \to \Theta^*$ is extended additively to $\mathscr{P}(\Sigma^*)$ by $\tau(L)= \bigcup_{f \in L} \tau(f)$ for every L in $\mathscr{P}(\Sigma^*)$.

Unless otherwise stated, we treat in this paper substitutions $\tau : \Sigma^* \to \Sigma^*$ and we call such $\tau$ a substitution on $\Sigma^*$. In this case we define for every integer n the products $\tau^{n+1}=\tau(\tau^n)$ to be those of relations ; these products are again substitutions. We shall use the following notations :

$\tau^*= \bigcup_{k \geq 0} \tau^k$ and $\tau^+=\tau(\tau^*)$, where $\tau^0$ is the identity mapping of $\Sigma^*$.

Let $u=x_1 x_2 \ldots x_\ell$ $x_i \in \Sigma$ and $v=s_1 s_2 \ldots s_\ell$ $s_i \in \Sigma^*$ be two words. The word v is said to be a descendant of u if v belongs to $\tau^n(u)$ for some positive integer n. The derivation $\delta$ from u to v is an $\ell$-tuple of pairs

$\delta = ((x_1,s_1),(x_2,s_2),\ldots,(x_\ell,s_\ell))$ where $s_i \in \tau^h(x_i)$ $i=1,2,\ldots,\ell$.

A substitution $\tau$ on $\Sigma^*$ is said to be _finite_ (resp. _rational_, _context free_) if for every a of $\Sigma$, $\tau(a)$ is a finite (resp. rational, context free) subset of $\Sigma^*$.

In the literature on L-system a pair $(\Sigma,\tau)$ is called a _0L-scheme_ if $\tau$ is a finite substitution on $\Sigma^*$.

Let alph be the function from $\Sigma^*$ into $\mathcal{F}(\Sigma)$ defined as follows : for any word f in $\Sigma^*$, alph(f) is the smallest subset S of $\Sigma$ such that f is in S*. The canonical additive extension of alph is thus a function from $\mathcal{F}(\Sigma^*)$ into $\mathcal{F}(\Sigma)$.

Let $\tau$ be a substitution on $\Sigma^*$. The _alphabetical projection_ of $\tau$, denoted by $\psi_\tau$, is the mapping from $\Sigma$ into $\mathcal{F}(\Sigma)$ defined by

$$\psi_\tau(a)=\text{alph}(\tau(a)).$$

The canonical additive extension of $\psi_\tau$ is a function from $\mathcal{F}(\Sigma)$ into itself and is again denoted by $\psi_\tau$.

LEMMA 1.1 [7] : $\psi_\tau^2 = \psi_{\tau^2}$ .

DEFINITION : A substitution $\tau$ on $\Sigma^*$ is said to be _alphabetically stable_ (or _stable_ for short) if the following two conditions hold :

    i) $\psi_\tau = \psi_\tau^2$

    ii) For every a in $\Sigma$, if 1 is in $\tau^+(a)$ then 1 is in $\tau(a)$.

PROPOSITION 1.2 [7] : _For every substitution_ $\tau$ _there exists an integer_ r _such that_ $\tau^r$ _is stable._

## 2. Repeatable words

DEFINITION : Let $\tau$ be a substitution on $\Sigma^*$. A word w of $\Sigma^*$ is <u>repeatable</u> for $\tau$ if it is a descendant of itself, i.e., $w \in \tau^+(w)$. We denote by $P(\tau)$ the set of repeatable words for $\tau$:

$$P(\tau) = \{w \mid w \in \tau^+(w)\}.$$

From the above definition a repeatable word u has at least one derivation $\delta$ from u to u. If $u = x_1 x_2 \ldots x_\ell$ and $\delta = ((x_1, s_1), (x_2, s_2), \ldots, (x_\ell, s_\ell))$ then u is factorized into $u = s_1 s_2 \ldots s_\ell$. We sometimes call this a factorization by $\delta$.

Let $R(\tau)$ be the set of recurrent words for substitution on $\Sigma^*$ [5,7]. That is, $R(\tau) = \{w \in \Sigma^* \mid w \in \tau^+(f) \text{ for any } f \in \tau^+(w)\}$. Then it is clear that $P(\tau) \supset R(\tau)$, and that a repeatable word is a very natural extension of that of recurrent words.

The followings are immediate consequences of the definition and have the corresponding version for recurrent words in [7], so we omit the proofs:

LEMMA 2.1 : $P(\tau)$ <u>is closed under product.</u>

LEMMA 2.2 : <u>For any positive integer</u> n, $P(\tau) = P(\tau^n)$.

## 3. Classification of letters

DEFINITION : Let $\tau$ be a substitution on $\Sigma^*$. A letter x of $\Sigma$ is said to be <u>vital</u> for $\tau$ if 1 is not a descendant of x,

i.e., if $1 \notin \tau^+(x)$. We denote by V the set of vital letters.
The set of non-vital letters is denoted by N, i.e., $N = \Sigma \setminus V$, or
equivalently, $N = \{x \mid 1 \in \tau^+(x)\}$.

PROPERTY 3.1 : Let $u = x_1 x_2 \ldots x_\ell$ $x_i \in \Sigma$ be repeatable word. Let
$\delta = ((x_1, s_1), (x_2, s_2), \ldots, (x_\ell, s_\ell))$ be a derivation from u to
u. Then

1. If $x_i$ is non-vital then $\left| s_i \right|_V = 0$.

2. If $x_i$ is vital then $x_i$ is the only vital letter contained
   in $s_i$.

DEFINITION : Let $\tau$ be a substitution on $\Sigma^*$. A letter x
in $\Sigma$ is said to be cyclic if there exist two words s and
t in N* such that sxt is in $\tau^+(x)$. The set of cyclic
letters is denoted by C.

DEFINITION : Let $\tau$ be a substitution on $\Sigma^*$. Let V be the
set of vital letters for $\tau$. $\tau$ is said to be vitality preserving
if for any x in $\Sigma$ and any u in $\tau(x)$, $|u|_V = |x|_V$.

Let $\tau$ be a substitution on $\Sigma^*$. Let V, N, and C the set
of vital, non-vital, and cyclic letters for $\tau$, respectively.
Consider the mapping $\tau'$ on (N∪C)* defined by

$$\begin{cases} \tau'(x) = \tau(x) \setminus \Sigma^* V \Sigma^* V \Sigma^* & \text{if } x \in C \cap V, \\ \tau'(x) = \tau(x) \setminus \Sigma^* V \Sigma^* & \text{if } x \in N. \end{cases}$$

It is easily seen that $\tau'$ is a well defined substitution on

(NUC)* and $\tau'$ is vitality preserving. We call $\tau'$ the vitality preserving substitution of $\tau$. The following is an immediate consequence of Property 3.1.

PROPERTY 3.2 : $P(\tau)=P(\tau')$.

That is, we may also assume, without loss of generality, that $\tau$ is vitality preserving to compute the set of repeatable words.

LEMMA 3.3 [7] : Let $\tau$ be stable substitution on $\Sigma^*$. For any f in $C^*$ there exists a word s such that f is a subword of s and s is in $\tau^k(f)$ for every positive integer k.

COROLLARY 3.4 : Let $\tau$ be a stable and vitality preserving substitution on $\Sigma^*$. For any x in C and any positive integer k

$$\tau^k(x) \subset \tau^{k+1}(x).$$

## 4. Letter position function

In this section we define a letter position function, which turns out to be very useful in the following discussion. We characterize the letters which appear in a repeatable word using the letter position function.

DEFINITION : Let $u=x_1 x_2 \ldots x_\ell$ $x_i \in \Sigma$ be a repeatable word for a substitution $\tau$ on $\Sigma^*$. Let $\delta=((x_1,s_1),(x_2,s_2),\ldots,(x_\ell,s_\ell))$ be the derivation from $u$ to $u$. The letter position function $\alpha : \{1,2,\ldots,\ell\} \rightarrow \{1,2,\ldots,\ell\}$ for $\delta$ is given by $\alpha(i)=j$ where $x_i$ is contained in $s_j$.

Informally, $\alpha$ indicates the ancestor of each letter of $u$ in the derivation $\delta$. It is very important, although obvious by the definition, that $\alpha$ is a non-decreasing function on the integer interval $[1,\ell]$. Then we have

PROPERTY 4.1 : Let $\alpha$ be a letter position function. $\alpha$ has at least one fixed point, i.e., there exists an $i$ in $\{1,2,\ldots,\ell\}$ such that $\alpha(i)=i$.

DEFINITION : Let $u=x_1 x_2 \ldots x_\ell$ be a repeatable word. Let $\alpha$ be the letter position function for a derivation $\delta$ from $u$ to $u$. Then $x_i$ is said to be repeatable in $u$ if $i$ is a fixed point of $\alpha$.

LEMMA 4.2 : Let $u=x_1 x_2 \ldots x_\ell$ be a repeatable word. If $x_i$ is vital then $x_i$ is repeatable in $u$.

Proof : Let $\alpha$ be the letter position function for a derivation $\delta$ from $u$ to $u$. Let $x_i$ be a vital letter and let $\alpha(i)=j$. From Property 3.1.1 $x_j$ must be vital. And then, from Property 3.1.2, we have $i=j$. $\square$

There is a very close relation between cyclic letter and repeatable letter in a repeatable word. Indeed we have

PROPOSITION 4.3 : If $x_i$ is repeatable in u then $x_i$ is in C.

Proof : Let $s_i$ be the segment produced by $x_i$ in a derivation from u to u. Then, from Lemma 4.2, $x_i$ is the only possible vital letter in $s_i$, even if $s_i$ contains any vital letters. Therefore, for some words f and g in N*, we have $s_i = fx_i g$ and $s_i$ is in $\tau^+(x_i)$. □

## 5. Characterization of repeatable words

DEFINITION : A repeatable word u for $\tau$ is said to be elementary if $u \neq 1$ and any factorization $u = u_1 u_2$ for repeatable words $u_1$ and $u_2$ implies $u_1 = 1$ or $u_2 = 1$. We denote by $P_1(\tau)$ the set of elementary repeatable words for $\tau$.

PROPERTY 5.1 : $P(\tau) = (P_1(\tau))^*$.

Proof : $P(\tau) \supset (P_1(\tau))^*$ is obvious by Lemma 2.1. $P(\tau) \subset (P_1(\tau))^*$ directly follows from the definition. □

LEMMA 5.2 : Let $u = x_1 x_2 \ldots x_\ell$ be a word in $P_1(\tau)$. There exists one and only one i such that $x_i$ is repeatable in u.

Proof : Assume there are exactly two cyclic letters $x_i, x_j$ $(i < j)$ for a letter position function $\alpha$. (We can show the lemma similarly in case there are more than two cyclic letters.)

As  i  and  j  are the fixed points of  $\alpha$,  there exists a positive integer  p  and an integer  k  $i \leqslant k < j$  such that

$$\alpha^p(n) = \begin{cases} i \text{ if } 1 \leqslant n \leqslant k \\ j \text{ if } k+1 \leqslant n \leqslant \ell. \end{cases}$$

That is,  $t_1 = x_1 x_2 \ldots x_k$  is in  $\tau^+(x_i)$  and  $t_2 = x_{k+1} \ldots x_\ell$  is in  $\tau^+(x_j)$.  Since  $u = t_1 t_2$  and  u  is in  $\tau^+(u)$,  $t_1$  is in  $\tau^+(t_1)$  and  $t_2$  is in  $\tau^+(t_2)$.  Thus  $t_1$  and  $t_2$  are repeatable.  This contradicts the fact that  u  is elementary.  ☐

The proof of the above lemma also shows

LEMMA 5.3 :  Let  $u = x_1 x_2 \ldots x_\ell$  be a word in  $P_1(\tau)$  and  $x_i$  be repeatable in  u.  Then  u  is in  $\tau^+(x_i)$.

PROPOSITION 5.4 :  Let  $\tau$  be a stable substitution on  $\Sigma^*$.  Let  C  and  N  be the set of cyclic and non-vital letters for  $\tau$, respectively.  Then

$$P_1(\tau) \subset \bigcup_{x \in C} ( \tau^+(x) \cap N^* x N^*) \subset P(\tau).$$

Proof : The left side inclusion is a corollary of the above Lemmas. Let  u  be a word in  $\tau^+(x) \cap N^* x N^*$  for some  $x \in C$.  Then  u  is written as  $u = sxt$  for some  st  in  $N^*$.  Since  $\tau$  is stable,  1  is in  $\tau(st)$,  and hence  u  is in  $\tau^+(sxt) = \tau^+(u)$.  ☐

COROLLARY 5.5 :  $P(\tau) = ( \bigcup_{x \in C} ( \tau^+(x) \cap N^* x N^*) )^*.$

Next we will give an effective finitary description of

$\bigcup_{x \in C} (\tau^+(x) \cap N^* x N^*)$. We need some additional definitions.

Throughout the remaining of this section let $\tau$ be a stable and vitality preserving substitution.

For any x in C let $C_\chi$ be the set defined by

$$C_\chi = \{y \mid \exists y_0, y_1, \ldots, y_n \in C \text{ such that } y_0 = y_n = x \ \exists k \ y_k = y \text{ and}$$
$$\text{for any } i \ 0 < i \leqslant n \ \exists s_i, t_i \in N^* \ s_i y_i t_i \in \tau(y_{i-1})\}.$$

Let $D_\chi, E_\chi, G_\chi,$ and $H_\chi$ be the sets defined by

$$E_\chi = \text{alph}(\{g \mid \exists f \in N^*, \ \exists y, z \in C_\chi \ fyg \in \tau(z)\})$$

$$H_\chi = \text{alph}(\{f \mid \exists g \in N^*, \ \exists y, z \in C_\chi \ fyg \in \tau(z)\})$$

$$D_\chi = \psi_\tau(E_\chi) \cap N \cap C$$

$$G_\chi = \psi_\tau(H_\chi) \cap N \cap C.$$


PROPOSITION 5.6 : <u>There exists an integer</u> n <u>such that for</u> <u>every</u> x <u>in</u> C

$$\tau^+(x) \cap N^* x N^* = \tau^n(G_\chi^* x D_\chi^*) \cap N^* x N^*.$$


In order to prove this Proposition we establish two lemmas, which can be shown quite similarly as in the proof of the corresponding results in [7].


LEMMA 5.7 : <u>Let</u> x <u>be a cyclic letter.</u> <u>For any word</u> fxg <u>in</u> $G_\chi^* x D_\chi^*$ <u>there exists a word</u> sxt <u>in</u> $\tau^+(x)$ <u>such that</u> fxg <u>is</u> <u>a subword of</u> sxt.


LEMMA 5.8 : <u>There exists an integer</u> k <u>with the following prop-</u> <u>erty :</u> Let x <u>be any cyclic letter,</u> sxt <u>any word in</u> $\tau^*(x) \cap N^* x N^*$,

_and_ c _any letter of_ s (resp. of t); _if_ c _is not in_ $\psi_\tau(G_\chi)$ (resp. in $\psi_\tau(D_\chi)$) _there exists a word_ uxv _in_ $\tau^k(x) \cap N^*xN^*$ _such that_ $sxt = s'uxvt'$ _and_ c _does not occur in_ s' (resp. in t').

Proof of Proposition 5.6 :

We first prove that

$$\tau^+(G_\chi^* x\ D_\chi^*) \subset \tau^+(x)$$

for any cyclic letter x. Let fxg be any word of $G_\chi^* x D_\chi^*$. From Lemma 5.7 there exist an integer k and a word sxt such that fxg is a subword of sxt and sxt is in $\tau^k(x)$. Now any descendant w of fxg is a descendant of sxt and thus of x.

Conversely let w be in $\tau^+(x) \cap N^*xN^*$ : there exist z in $C_\chi$, and s,t in N* such that w is in $\tau(szt)$, i.e., there exist two subwords $f = a_1 a_2 \ldots a_p$ and $g = b_1 b_2 \ldots b_q$ of s and t respectively and a factorization $w = s_1 s_2 \ldots s_p uyv t_1 t_2 \ldots t_q$ such that $s_i$ is in $\tau(a_i)$, $t_j$ in $\tau(b_j)$ for every i and j, and uyv is in $\tau(z)$. If all the $a_i$'s and the $b_j$'s are in $\psi_\tau(G_\chi)$ and $\psi_\tau(D_\chi)$ respectively w is in $\tau^2(G_\chi^* x D_\chi^*)$. If it is not the case it follows from Lemma 5.8 that there exists a factorization $szt = s'u'zv't'$ such that all the $a_i$ occurring in s' are in $\psi_\tau(G_\chi)$, all the $b_j$ occurring in t' are in $\psi_\tau(D_\chi)$ and uzv is in $\tau^k(z)$. Then w is in $\tau^{k+1}(G_\chi^* x D_\chi^*)$. $\square$

THEOREM 5.9 : The set of repeatable words for a rational (resp.context free) substitution is a rational (resp.context free) set.

Proof : The families of rational languages and context free languages are both closed under union, intersection with rational set, Kleene closure, and substitution. Then the theorem follows immedeately from Corollary 5.5 and Proposition 5.6. □

REFERENCES

[1] Herman, G.T. & Rozenberg, G. : Developmental systems and languages, North-Holland, Amsterdam, 1974.

[2] Herman, G.T. & Walker, A. : Context free languages in biological systems, Inter. J. of Computer Math., 4, 369-391, 1975.

[3] Král, J. : A modification of a substitution theorem and some necessary and sufficient conditions for sets to be context free, Math. Systems Teory, 4, 129-239, 1970.

[4] Lindenmayer, A. : Developmental systems without cellular interactions, their languages and grammars, J. Theor. Biol., 30, 455-484, 1971.

[5] Nishida, T. & Kobuchi, Y. : Recurrent strings in a 0L language, 数理解析研究所講究録, 421, 121 - 133, 1781.

[6] Rozenberg, G. & Salomaa, A. : The mathematical theory of L systems, Academic Press, New York, 1980.

[7] Sakarovitch, J., Nishida, T. & Kobuchi, Y. : Recurrent words for substitution, Technical report of Research Institute for Mathematical Sciences Kyoto University, RIMS-396, 1982.