

JANUSZ A. BRZozowski

Hierarchies of aperiodic languages

*Revue française d'automatique informatique recherche opérationnelle.
Informatique théorique*, tome 10, n° R2 (1976), p. 33-49

<http://www.numdam.org/item?id=ITA_1976__10_2_33_0>

© AFCET, 1976, tous droits réservés.

L'accès aux archives de la revue « Revue française d'automatique informatique recherche opérationnelle. Informatique théorique » implique l'accord avec les conditions générales d'utilisation (<http://www.numdam.org/conditions>). Toute utilisation commerciale ou impression systématique est constitutive d'une infraction pénale. Toute copie ou impression de ce fichier doit contenir la présente mention de copyright.

NUMDAM

Article numérisé dans le cadre du programme
Numérisation de documents anciens mathématiques

<http://www.numdam.org/>

HIERARCHIES OF APERIODIC LANGUAGES (*)

by JANUSZ A. BRZOZOWSKI ⁽¹⁾

Communiqué par J. F. PERROT

Abstract. — In recent years, considerable attention has been given to the family of aperiodic languages, also known as star-free languages and noncounting regular languages. Several interesting subfamilies of aperiodic languages have been studied and characterized by the properties of the corresponding syntactic semigroups. The study of such families has been systematized by examining the position of each family in certain natural hierarchies. This paper gives a brief survey of results in this area.

ACKNOWLEDGMENT

This work was done while the author at the University of Paris VI and VII, under the scientific exchange program between Canada and France. The financial assistance provided by this program is gratefully acknowledged.

Preliminary versions of this paper were given at the Séminaire P. Dubreil: Algèbre (28^e année, 1974/1975, n° 19, 11 p.) and at the École de Printemps d'Informatique théorique 1975.

1. INTRODUCTION

Before considering the family of aperiodic languages (also known as star-free languages), we recall some basic notions about the more general family of regular languages (also known as rational languages). Precise mathematical definitions of the concepts mentioned in this section will be postponed. Kleene [12] related regular languages to finite automata by showing that for each regular language there is a finite automaton recognizing it, and that each language recognized by a finite automaton is regular. This correspondence between regular languages and finite automata has been extended in two directions: The first of these is practical, relating finite automata to certain types of sequential circuits; the second extension is mathematical, relating regular languages to finite semigroups. Thus the interest in finite automata and regular languages is shared by the theories of sequential circuits, formal languages, computing machine models and finite semigroups.

(*) Reçu octobre 1975.

(¹) Computer Science Department University of Waterloo, Waterloo, Ontario, Canada.

The family of aperiodic languages is a proper subfamily of the family of regular languages, and is of interest from several points of view. First, it corresponds to feedback-free sequential circuits constructed with gates and set-reset flip-flops. Second, the corresponding automata are precisely the “permutation-free” automata. Third, the aperiodic languages constitute the subfamily of regular languages defined without the use of the star operator (hence “star-free”). Fourth, the corresponding semigroups are precisely the “group-free” semigroups. The basic relationship between star-free languages and group-free semigroups was shown by Schützenberger [19, 20]. Several other interesting characterizations are investigated in the monograph by McNaughton and Papert [14]; however, the above-mentioned four points of view provide ample motivation for the study of aperiodic languages.

Within the family of aperiodic languages, a number of interesting subfamilies have been studied. Among the earliest considered are the definite languages [2, 12, 17], introduced by Kleene. These are characterized by the property that the membership of a word in a definite language depends only on the length- k suffix of that word, for some integer $k \geq 0$. This family is significant from the sequential circuit point of view, since definite languages correspond to feedback-free circuits constructed with gates and unit delays. Also there exist elegant (though somewhat more technical) characterizations of definite automata [17], and of the corresponding semigroups (discussed later).

Left-right duality led naturally to reverse definite languages [2], where membership is determined by prefixes rather than by suffixes. The idea of testing both the prefix and suffix of a word, thus obtaining the “generalized definite” languages, is due to Ginzburg [11]. All these languages are special cases of the locally-testable languages first studied by McNaughton and Papert [14]. These are of independent interest and have appeared previously in formal language theory [7, 15]. (Thus, for example, each context-free language is a homomorphic image of a Dyck language and a certain locally testable language.)

A systematic study of subfamilies of aperiodic languages was begun by Cohen and Brzowski [8] and continued by Brzowski and Simon [6]. The first step consisted of the introduction of “dot-depth” as a measure of complexity of aperiodic languages. The following motivation led to these concepts. Feedback-free networks of gates, i.e. combinational circuits, constitute the simplest and degenerate forms of sequential circuits. Combinational networks are, of course, characterized by Boolean functions. This suggested that (a) all Boolean operations should be considered together when studying the formation of aperiodic languages from the letters of the alphabet, and (b) since concatenation (or “dot” operator) is linked to the sequential rather than the combinational nature of a language, the number of conca-

tenation levels required to express a given aperiodic language should provide a useful measure of complexity. (Such reasoning is not precise, and is given here only as an intuitive guide. However, subsequent studies provided considerable evidence that this is indeed a useful approach.) As it turns out, locally testable languages require only one level of concatenation, i. e. are of "dot-depth" one. A finer measure of complexity is obtained when one also takes into account the number n of factors used, considering concatenation as an n -ary operation [6]. With this refinement the families of definite, reverse definite, generalized definite and locally testable languages appear naturally in a hierarchy of families whose union is the family of languages of dot-depth one. Moreover, the families of semigroups corresponding to these families of languages appear in a rather natural hierarchy of semigroups, as we shall see later.

Simon [23, 24] showed the correspondence between certain languages of depth-one and \mathcal{S} -trivial monoids, thus providing a link to classical semigroup theory. This is pursued further in [4, 21].

There remains a large number of open problems, and this paper has been written the hope that it will stimulate further work in this area. The proofs of several key results stated here are quite involved and lengthy. We do not repeat them here, since it is our aim to provide a brief overview of the subject, the main results, and the open problems.

References to specific results are given in the text. For further general background on aperiodic languages see the papers [1, 5, 16, 19, 20] and the books by Eilenberg [9] and McNaughton and Papert [14].

2. NOTATION

If A is a finite, non-empty alphabet, A^+ (respectively A^*) is the free semigroup (respectively free monoid) generated by A . The empty word is denoted by 1 , and Φ is the empty set. Any subset L of A^* is a language. The length of a word $w \in A^*$ is denoted by $|w|$. The cardinality of a set X is denoted $\text{card } X$. The symbol $\stackrel{\Delta}{=}$ means "is by definition".

Given languages $L, L' \subset A^*$, the following are also languages:

$$L \cup L' \quad (\text{union}),$$

$$L \cap L' \quad (\text{intersection}),$$

$$L = A^* - L \quad (\text{complement}),$$

$$L.L' \stackrel{\Delta}{=} \{w; w = uu', u \in L, u' \in L'\} \quad (\text{concatenation or product}),$$

$$L^+ \stackrel{\Delta}{=} \bigcup_{n \geq 1} L^n \quad (\text{the subsemigroup of } A^* \text{ generated by } L)$$

and

$$L^* \stackrel{\Delta}{=} \bigcup_{n \geq 0} L^n = L^+ \cup 1 \quad (\text{the submonoid of } A^* \text{ generated by } L).$$

Let \mathcal{U}_A (or simply \mathcal{U} when A is understood) be the family of all languages over A . Evidently, \mathcal{U}_A is a Boolean algebra under union, intersection and complement, and a monoid under concatenation.

Let $\mathcal{L}_A \stackrel{\Delta}{=} \{ \{ a \}; a \in A \}$ and let $\mathcal{W}_A \stackrel{\Delta}{=} \{ \{ w \}; w \in A^* \}$. Let \mathcal{F}_A be the family of all finite languages, and $\mathcal{C}_A \stackrel{\Delta}{=} \{ L \subset A^*; \bar{L} \in \mathcal{F}_A \}$ the family of cofinite languages.

For a given family \mathcal{X} of languages, consider the following properties:

- (a) $(L, L' \in \mathcal{X}) \Rightarrow (L \cup L' \in \mathcal{X})$;
- (b) $(L \in \mathcal{X}) \Rightarrow (\bar{L} \in \mathcal{X})$;
- (c) (i) $\{ 1 \} \in \mathcal{X}$,
- (ii) $L, L' \in \mathcal{X} \Rightarrow LL' \in \mathcal{X}$;
- (d) $(L \in \mathcal{X}) \Rightarrow (L^* \in \mathcal{X})$.

It is well known [9, 12] that the family of regular or rational languages can be defined as the smallest family containing \mathcal{L}_A and satisfying (a), (c) and (d), and that this family also satisfies (b).

Aperiodic languages can be defined as the smallest family containing \mathcal{L}_A , and satisfying (a), (b) and (c).

As we have said before, in the study of aperiodic languages it is useful to separate the closure under Boolean operations from the closure under concatenation. For any family $\mathcal{X} \subset \mathcal{U}$, denote by $\mathcal{X}B$ the Boolean algebra generated by \mathcal{X} , i. e. the smallest family containing \mathcal{X} and satisfying (a) and (b). Similarly, $\mathcal{X}M$ denotes the monoid generated by \mathcal{X} , i. e. the smallest family containing \mathcal{X} and satisfying (c).

3. APERIODIC LANGUAGES OVER A ONE-LETTER ALPHABET

For $A = \{ a \}$, the family \mathcal{A}_a of aperiodic languages is particularly simple. (We use \mathcal{A}_a for $\mathcal{A}_{\{a\}}$, etc.) We have $\mathcal{L}_a M = \mathcal{W}_a = \{ \{ a^n \}; n \geq 0 \}$. Define

$$\mathcal{B}_a \stackrel{\Delta}{=} \mathcal{L}_a M B.$$

Since \mathcal{B}_a must be closed under union, it contains all finite languages, $\mathcal{B}_a \supset \mathcal{F}_a$. Closure under complementation implies $\mathcal{B}_a \supset \mathcal{C}_a$. One verifies that each cofinite language L can be written in the form

$$L = F \cup a^n a^*,$$

for some $n \geq 0$ and some $F \in \mathcal{F}_a$. It is now clear that $\mathcal{F}_a \cup \mathcal{C}_a$ is closed under both complementation and union, i. e. that it is a Boolean algebra. Thus we conclude that

$$\mathcal{B}_a = \mathcal{F}_a \cup \mathcal{C}_a.$$

Moreover, note that concatenation of languages over a one-letter alphabet is commutative. Using this and the form (1) for cofinite languages, one verifies that \mathcal{B}_a is closed under concatenation also, i. e.:

$$\mathcal{B}_a M = \mathcal{B}_a.$$

This implies that all aperiodic languages are in \mathcal{B}_a , i. e.:

$$\mathcal{A}_a = \mathcal{B}_a.$$

Thus a language over a one-letter alphabet is aperiodic if, and only if, it is either finite or cofinite.

If we start by closing \mathcal{L}_a under Boolean operations first, we find

$$\tilde{\mathcal{B}}_a \stackrel{\Delta}{=} \mathcal{L}_a B = \{\Phi, \{a\}, a^*, a^* - a\}.$$

Next note that

$$\tilde{\mathcal{M}}_a \stackrel{\Delta}{=} \tilde{\mathcal{B}}_a M \supset \{\{a\}, a^*\} M = (\mathcal{L}_a \cup a^*) M = (\mathcal{W}_a \cup a^*) M$$

and

$$\mathcal{L}_a B M B = \mathcal{M}_a B \supset \mathcal{F}_a \cup \tilde{\mathcal{C}}_a = \mathcal{B}_a = \mathcal{A}_a$$

because each finite and cofinite language can be expressed as an element of $(\mathcal{W}_a \cup a^*) MB$. Since, obviously,

$$\tilde{\mathcal{M}}_a B \subset \mathcal{A}_a,$$

we have

$$\mathcal{A}_a = \mathcal{B}_a = \tilde{\mathcal{M}}_a B.$$

These observations are summarized in Figure 1. For each inclusion, we provide an example of a language which proves the inclusion is proper.

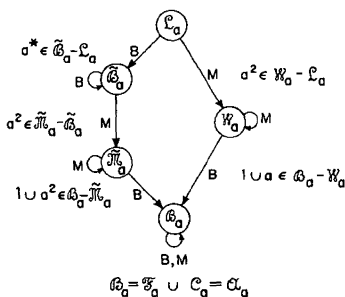


Figure 1.
Aperiodic languages over a one-letter alphabet.

4. INITIAL PHENOMENA [8]

We now assume that the alphabet A is fixed, and $\text{card } A > 1$. We use \mathcal{L} for \mathcal{L}_A , etc. As in the one-letter case, we have $\mathcal{W} \stackrel{\Delta}{=} \mathcal{L} M$ and

$$\mathcal{B}_0 \stackrel{\Delta}{=} \mathcal{L} M B = \mathcal{F} \cup \mathcal{C}.$$

However, $\mathcal{B}_0 M \neq \mathcal{B}_0$ since (for $A = \{a, b\}$) the language $\{a, b\}^* a = \bar{\Phi} \cdot a$ is in \mathcal{B}_0 , but is neither finite nor cofinite. Thus we proceed to define

$$\mathcal{M}_1 \stackrel{\Delta}{=} \mathcal{B}_0 M \quad \text{and} \quad \mathcal{B}_1 \stackrel{\Delta}{=} \mathcal{M}_1 B.$$

We will return to these families later. For now observe that

$$\mathcal{B}_1 = (\mathcal{F} \cup \mathcal{C}) M B = (\mathcal{W} \cup A^*) M B,$$

since each cofinite language can be written $L = F \cup A^n A^*$, for some $n \geq 0$ and $F \in \mathcal{F}$. Hence L can be written as a union of products where each factor is either A^* or it is in \mathcal{W} .

If we close \mathcal{L} under Boolean operations first, we find

$$\tilde{\mathcal{B}}_0 \stackrel{\Delta}{=} \{L; L \subset A\} \cup \{L; \bar{L} \subset A\}.$$

Thus $\tilde{\mathcal{B}}_0$ is a finite Boolean algebra with $\mathcal{L} \cup \bar{A}$ as the set of atoms. Note that $\tilde{\mathcal{B}}_0 \supset \mathcal{L} \cup A^*$. Next

$$\tilde{\mathcal{M}}_1 \stackrel{\Delta}{=} \tilde{\mathcal{B}}_0 M \supset (\mathcal{L} \cup A^*) M = (\mathcal{W} \cup A^*) M.$$

Thus $\tilde{\mathcal{M}}_1 B \supset (\mathcal{W} \cup A^*) M B = \mathcal{B}_1$. Conversely,

$$\mathcal{B}_1 = \mathcal{L} M B M B \supset \mathcal{L} B M B = \mathcal{M}_1 B,$$

and

$$\mathcal{B}_1 = \tilde{\mathcal{M}}_1 B.$$

These properties are summarized in Figure 2. It is seen that, except for the few initial differences, it is not important whether \mathcal{L} is closed under B or M first, since the two sequences coincide from \mathcal{B}_1 on.

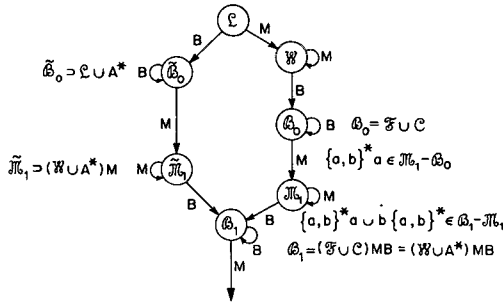


Figure 2. Initial families for card $A \geq 2$.

5. THE DOT-DEPTH HIERARCHY [8]

The sequence (\mathcal{B}_i) of Boolean algebras, defined below, is called the *dot-depth hierarchy*. Let

$$\mathcal{B}_0 \stackrel{\Delta}{=} \mathcal{L} M B,$$

$$\mathcal{B}_{n+1} \stackrel{\Delta}{=} \mathcal{B}_n M B = \mathcal{B}_0 (M B)^n = \mathcal{L} (M B)^{n+1} \quad \text{for } n \geq 0.$$

For each aperiodic language L , there exists $n \geq 0$ such that $L \in \mathcal{B}_n$; hence

$$\mathcal{A} = \bigcup_{n \geq 0} \mathcal{B}_n.$$

The “position” of a language in the dot-depth hierarchy can be used as a measure of its complexity. Define the *dot-depth* (or simply *depth*) of a language L by

$$d(L) = 0 \quad \text{if } L \in \mathcal{B}_0,$$

$$d(L) = n \quad \text{if } L \in \mathcal{B}_n - \mathcal{B}_{n-1} \quad \text{for } n > 0.$$

The depth $d(L)$ corresponds to the minimum number of concatenation levels that must be used to generate L from languages in \mathcal{B}_0 . Also, $\tilde{\mathcal{B}}_0$ can be used instead of \mathcal{B}_0 since $\tilde{\mathcal{B}}_0 M B = \mathcal{B}_0 M B$; however, \mathcal{B}_0 appears to be a more natural starting point (see Fig. 3).

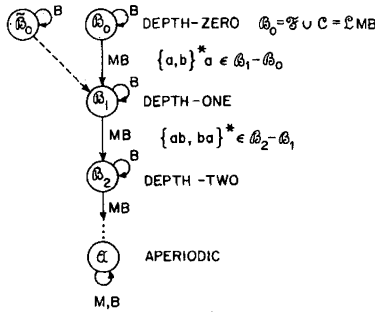


Figure 3.
The dot-depth hierarchy.

It has been shown recently by Brzozowski and Knast (*) that the dot-depth hierarchy is infinite for card $A > 1$.

For $A = \{ a, b, c \}$, the language

$$L_1 = \{ a, b, c \}^* b a^* = \overline{\Phi} . b . \overline{\overline{\Phi}} . \{ b, c \} . \overline{\overline{\Phi}}$$

is of depth 2, i. e. $L_1 \in \mathcal{B}_2 - \mathcal{B}_1$ [8]. An example over a two letter alphabet is $L_2 = \{ ab, ba \}^*$.

An upper bound for $d(L)$ has been found as follows [8]. Let n be the number of states in the reduced deterministic finite automaton \mathcal{U}_L recognizing L . Let $i_a(L)$ be the number of distinct states in input column a , $a \in A$, of the state table of \mathcal{U}_L . Further, let

$$i(L) = \max \{ i_a(L); a \in A \text{ and } i_n(L) \neq n \}.$$

Then $d(L) \leq i(L) + 1$.

This bound is met by L_1 above. On the other hand, let $L_n = a^n a^*$, $n \geq 0$ be over $A = \{ a \}$. One verifies that $i(L_n) = n - 1$, although L_n is cofinite, and $d(L_n) = 0$.

6. THE DEPTH-ONE FINITE-COFINITE HIERARCHY [6]

In the dot-depth hierarchy, $\mathcal{B}_1 = \mathcal{B}_0 MB$, i. e. a language in \mathcal{B}_1 is a Boolean function of products of any number of factors from \mathcal{B}_0 . Thus only one level of concatenation is required, but this concatenation is unlimited in the number of

(*) J. A. BRZOZOWSKI and R. KNAST, *The Dot-Depth Hierarchy of Star-Free Languages is Infinite*, Research Report CS-76-23, Computer Science Dept., University of Waterloo, Waterloo, Ont., Canada; April, 1976.

factors. A finer measure of complexity is obtained by limiting the number of factors as follows. Let

$$\beta_n \stackrel{\Delta}{=} \mathcal{B}_0^n B \quad \text{for } n \geq 1.$$

Then

$$\beta_n \subset \beta_{n+1} \quad \text{and} \quad \mathcal{B}_1 = \bigcup_{n \geq 1} \beta_n.$$

A number of subfamilies of aperiodic languages that have been studied appear naturally in the sequence $\mathcal{B}_0 = \beta_1 \subset \beta_2 \subset \dots \subset \mathcal{B}_1$ which we refer to as the *(depth-one) finite-cofinite hierarchy*. We will also need:

$$\begin{aligned} \beta_{2L} &\stackrel{\Delta}{=} (\mathcal{F}^2 \cup \mathcal{C}\mathcal{F} \cup \mathcal{C}^2) B = (\mathcal{F} \cup \mathcal{C}\mathcal{F} \cup \mathcal{C}) B \subset \beta_2, \\ \beta_{2R} &\stackrel{\Delta}{=} (\mathcal{F}^2 \cup \mathcal{F}\mathcal{C} \cup \mathcal{C}^2) B = (\mathcal{F} \cup \mathcal{F}\mathcal{C} \cup \mathcal{C}) B \subset \beta_2. \end{aligned}$$

An alternate description of the β families is the following:

$$\begin{aligned} \beta_1 &= \mathcal{B}_0 = \mathcal{W} B = (\mathcal{W} \cup A^*) B, \\ \beta_{2L} &= (\mathcal{W} \cup A^* \mathcal{W}) B, \\ \beta_{2R} &= (\mathcal{W} \cup \mathcal{W} A^*) B \\ \beta_n &= (\mathcal{W} \cup A^*)^n B = (\mathcal{F} \cup \mathcal{C})^n B \quad \text{for } n \geq 1, \end{aligned}$$

where $A^* \mathcal{W} = \{ A^* L; L \in \mathcal{W} \}$, etc. These claims are easily verified. One can also show [6] that

$$(\mathcal{W} \cup A^*)^{2n+1} B \supset (\mathcal{W} \cup A^*)^{2n+2} B \quad \text{for } n \geq 1.$$

Therefore, $\beta_{2n+2} = \beta_{2n+1}$; however, $\beta_{2n+3} \neq \beta_{2n+1}$ for all $n \geq 1$ [23].

A language is *definite* [2, 12, 17] if, and only if, is in β_{2L} , *reverse definite* [2, 11] if, and only if, it is in β_{2R} , *generalized definite* [11, 22] if, and only if, it is in β_2 , and *locally testable* [14] if, and only if, it is in β_3 . The original definitions of these families of languages were somewhat different; however, the equivalence of the definitions is easily proved [6], and the present formulation appears more natural. We reconsider these families later.

The statements about the finite-cofinite hierarchy are summarized in Figure 4.

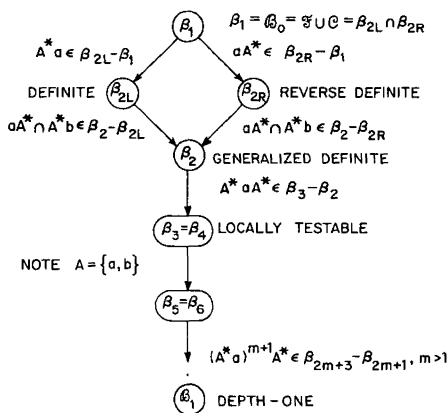


Figure 4. The depth-one finite-cofinite hierarchy.

7. THE γ_1 HIERARCHY [23, 24]

The languages introduced here play a key role in the family of depth-one languages. We introduce a family $\alpha_{1,1}$ of languages (the reason for this notation is explained in Section 9) such that, if $L \in \alpha_{1,1}$, the membership of a word x in L can be determined solely by the set of letters appearing in x . Define

$$x\alpha \stackrel{\Delta}{=} \{ a \in A; x = uav \text{ and } u, v \in A^* \}$$

to be the “alphabet” of $x \in A^*$.

For $x, y \in A^*$, let $x \equiv_\alpha y$ if, and only if, $x\alpha = y\alpha$. The relation \equiv_α is a congruence of finite index on A^* , there being one congruence class $([x]_\alpha)$ for each subset of A . We have

$$[x]_\alpha = \left(\bigcap_{a \in x\alpha} A^* a A^* \right) \cap \left(\bigcap_{a \notin x\alpha} \overline{A^* a A^*} \right).$$

Now define $\alpha_{1,1} \stackrel{\Delta}{=} \{ [x]_\alpha; x \in A^* \} B$, and let $A^* \mathcal{L} A^* \stackrel{\Delta}{=} \{ A^* a A^*; a \in A \}$. One verifies that $\alpha_{1,1} = (A^* \mathcal{L} A^*) B$. For technical reasons, we use the family $\mathcal{G} \stackrel{\Delta}{=} A^* \cup A^* \mathcal{L} A^*$ as a generating set for $\alpha_{1,1}$. Note that $\mathcal{G}^m \subset \mathcal{G}^{m+1}$ for $m \geq 1$, and we will use the convention $\mathcal{G}^0 = \{ \Phi \}$. Let $\alpha_{m,1} \stackrel{\Delta}{=} \mathcal{G}^m B$ and $\gamma_1 \stackrel{\Delta}{=} \mathcal{G} M B$. We find $\gamma_1 = \mathcal{G} M B = \bigcup_{m \geq 0} \mathcal{G}^m B = \bigcup_{m \geq 0} \alpha_{m,1}$ and the sequence $\alpha_{0,1} \subset \alpha_{1,1} \subset \alpha_{2,1} \subset \dots \subset \gamma_1$ will be called the γ_1 hierarchy.

An alternate description of $\alpha_{m,1}$ is obtained by using the “shuffle” operator \sqcup [9, 10]. For $w = a_1 a_2 \dots a_m \in A^*$,

$$A^* \sqcup w = A^* \sqcup (a_1 a_2 \dots a_m) \stackrel{\Delta}{=} A^* a_1 A^* a_2 A^* \dots a_m A^*$$

and

$$A^* \sqcup \mathcal{W} = \{ A^* \sqcup w; w \in A^* \}.$$

Let $\mathcal{W}_{\leq m} = \{ \{ w \}; w \in A^* \text{ and } |w| \leq m \}$. One verifies that

$$\alpha_{m,1} = (A^* \sqcup \mathcal{W}_{\leq m}) B, \text{ for } m \geq 0, \text{ and } \gamma_1 = (A^* \sqcup \mathcal{W}) B.$$

Over a two-letter alphabet $A = \{ a, b \}$, the γ_1 languages can be viewed as a generalization of finite-cofinite languages [3]. Let

$$\mathcal{L}_{\oplus} \stackrel{\Delta}{=} \{ \{ a \}, \{ b \}, a^+, b^+ \}$$

and $\mathcal{W}_{\oplus} \stackrel{\Delta}{=} \mathcal{L}_{\oplus} M$ be the generalization of \mathcal{L} and \mathcal{W} , respectively. Let \mathcal{F}_{\oplus} be the closure of the family \mathcal{W}_{\oplus} under finite unions and let $\mathcal{C}_{\oplus} \stackrel{\Delta}{=} \{ L; \bar{L} \in \mathcal{F}_{\oplus} \}$. Then it can be shown that

$$\gamma_1 = \mathcal{F}_{\oplus} \cup \mathcal{C}_{\oplus} = \mathcal{L}_{\oplus} M B.$$

Furthermore, the initial phenomena of Figure 2 have their counterpart here, for

$$\mathcal{L}_{\oplus} B M B = \mathcal{L}_{\oplus} M B M B.$$

8. THE LOCALLY-TESTABLE HIERARCHY [6, 23]

It can be shown that the membership of a word x in a locally-testable language L is determined solely by the first $k-1$ letters of x , the last $k-1$ letters of x and the set of words of length k that appear in x . Formally, $f_k(x)$ [respectively $t_k(x)$] is x , if $|x| \leq k$, and it is the prefix (respectively suffix) of x of length k otherwise. Let

$$m_k(x) \stackrel{\Delta}{=} \{ w \in A^*; x = uwv \text{ and } |w| = k \}.$$

For $x, y \in A^*$ and $k > 0$, define the congruence

(*) $x \sim_k y$ if and only if,

$$f_{k-1}(x) = f_{k-1}(y), t_{k-1}(x) = t_{k-1}(y) \quad \text{and} \quad m_k(x) = m_k(y).$$

If $[x]_k$ is the congruence class containing x , let $\alpha_{1,k} \stackrel{\Delta}{=} \{ [x]_k; x \in A^* \} B$, be the family of k -testable languages. The reason for this notation will soon be explained. Note however that it is consistent with that of Section 7, because L is a 1-testable if, and only if, $L \in \alpha_{1,1}$. One verifies that $\alpha_{1,k} \subset \alpha_{1,k+1}$ and that $\beta_3 = \bigcup_{k \geq 1} \alpha_{1,k}$.

If, in the definition (*) of \sim_k , we remove the condition $m_k(x) = m_k(y)$, we obtain the family of *k-generalized-definite* languages which we denote by $\alpha_{0,k}$. One verifies that $\beta_2 = \bigcup_{k \geq 1} \alpha_{0,k}$. If only t_{k-1} is tested we obtain the family of definite languages, and reverse definite languages are obtained by testing f_{k-1} . See Figure 4 for the location of these languages in the depth-one finite-cofinite hierarchy.

9. SIMON'S DEPTH-ONE HIERARCHY [23]

The membership of a word x in a language L of depth 1 can be determined by testing $f_{k-1}(x)$, $t_{k-1}(x)$ and the set $\mu_{m,k}$ of m -tuples of words of length k that "occur" in x . Thus depth-one languages are generalizations of both the k -testable and $\alpha_{m,1}$ languages; the locally testable and γ_1 hierarchies turn out to be "orthogonal".

More formally, let $W = (w_1, \dots, w_m)$ be an m -tuple of words of length k . We say that W occurs in x if, and only if, there exist words $u_1, \dots, u_m, v_1, \dots, v_m$ such that $|u_1| < |u_2| < \dots < |u_m|$ and $x = u_i w_i v_i$, for $i = 1, \dots, m$. Let

$$\begin{aligned} \mu_{m,k}(x) &= \{ W \mid W = (w_1, \dots, w_m), |w_1| = \dots \\ &= |w_m| = k \text{ and } W \text{ occurs in } X \}. \end{aligned}$$

By convention $\mu_{0,k} = \Phi$ for all $k \geq 1$. Note that $\mu_{m,k}(x) = \Phi$ if, and only if, $|x| < m+k-1$. For $x, y \in A^*$, $m \geq 0$, $k \geq 1$ define $x \sim_k y$ if, and only if,

$$(a) \quad x = y \quad \text{if} \quad |x| < m+k-1$$

or

$$(b) \quad f_{k-1}(x) = f_{k-1}(y), \quad t_{k-1}(x) = t_{k-1}(y)$$

and

$$\mu_{m,k}(x) = \mu_{m,k}(y), \quad \text{otherwise.}$$

The relation \sim_k is a congruence of finite index on A^* . Let

$$\alpha_{m,k} \stackrel{\Delta}{=} \{ {}_m[x]_k; x \in A^* \} B.$$

One verifies that this is consistent with the previous definitions.

The hierarchy defined by \sim_k is illustrated in Figure 5, where $\gamma_i = \bigcup_{m \geq 0} \alpha_{m,k}$ and (one verifies that) $\beta_{2m+1} = \bigcup_{k \geq 1} \alpha_{m,k}$ for $m \geq 1$ (the case of β_2 is somewhat degenerate). All the hierarchies shown are known to be infinite.

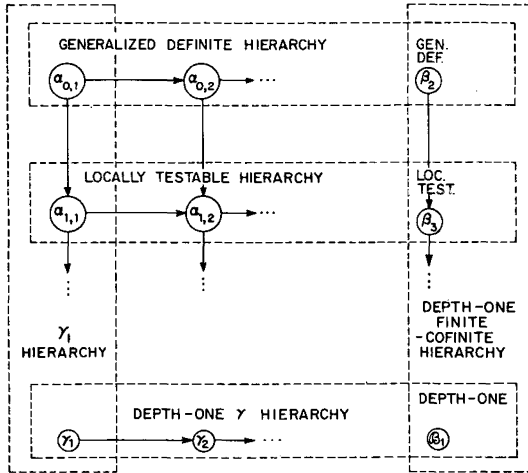


Figure 5.
Simon's depth-one hierarchy.

10. SYNTACTIC SEMIGROUPS AND MONOIDS

The congruence $m \sim_k$ of the previous section can be viewed as a characterization of the family \mathcal{B}_1 of depth-one languages since $L \in \mathcal{B}_1$ if, and only if, there exist $m \geq 0$ and $k \geq 1$ such that L is a union of congruence classes of $m \sim_k$. However, the problem is to decide effectively, given a regular language L , whether such m and k exist. In certain cases described below a decision procedure is available through a characterization of the syntactic semigroup or monoid of L .

For $L \subset A^*$ the syntactic congruence, \equiv_L , is defined by $x \equiv_L y$, if, and only if, for all $u, v \in A^*$, $(uxv \in L) \Leftrightarrow (uyv \in L)$. The quotient monoid A^*/\equiv_L is called the syntactic monoid M_L of L , and A^+/\equiv_L is the syntactic semigroup, S_L .

At first the difference between S_L and M_L appears to be rather trivial; however, it is essential in some cases to distinguish between M_L and S_L .

It is well-known that L is regular if, and only if, M_L is finite. It has been shown by Schützenberger [19, 20] that L is aperiodic if, and only if, M_L is finite and group-free (contains no groups other than the trivial one-element groups).

A number of families of languages in Simon's hierarchy have been characterized by the properties of their syntactic monoids. In this connection, the γ_1 hierarchy plays a key role. The following is known [23] :

- (1) $L \in \alpha_{0,1} = \{ \emptyset, A^* \}$ if, and only if, $M_L = 1$;

- (2) $L \in \alpha_{1,1}$ if, and only if, M_L is finite and idempotent and commutative;
- (3) $L \in \gamma_1$ if, and only if, M_L is finite and \mathcal{S} -trivial, i. e. for all $m, m' \in M_L$, $(M_L m M_L = M_L m' M_L)$ implies $(m = m')$.

These properties appear to carry over to the depth-one finite-cofinite hierarchy as follows:

“ $L \in \alpha_{m,1}$ if, and only if, M_L has property P ” seems to correspond to “ $L \in \beta_{2m+1}$ if, and only if, for each idempotent $e \in S_L$ the submonoid $e S_L e$ has property P ”.

The following evidence supports this statement:

- (1*) $L \in \beta_2$ if, and only if, S_L is finite and $e S_L e = e$ [6, 18, 25].
- (2*) $L \in \beta_3$ if, and only if, S_L is finite and $e S_L e$ is idempotent and commutative [6, 13, 25, 26, 27].
- (3*) If $L \in \mathcal{B}_1$, then S_L is finite and $e S_L e$ is \mathcal{S} -trivial [23].

As can be seen, the results are quite fragmentary, and the proofs of these results are quite complex. This approach appears to be very fruitful not only for classifying languages, but also monoids.

Simon’s β hierarchy begins with generalized definite languages. Finite-cofinite, definite and reverse definite languages represent special cases, and can also be characterized by the corresponding semigroups [10, 27]. This is summarized in Figure 6.

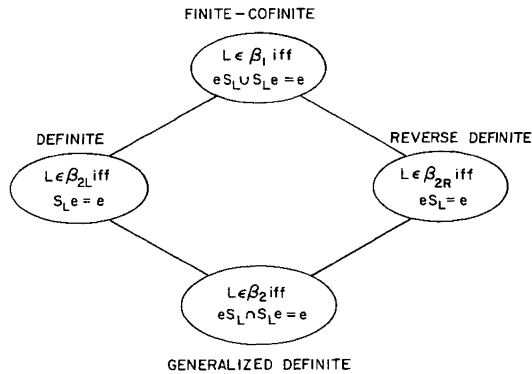


Figure 6.
Characterization of generalized definite languages.

One can generalize these ideas as follows [4]. For any monoid M and $m \in M$ define $P_m \stackrel{\Delta}{=} \{ m'; m \in M m' M \}$ and $M_m \stackrel{\Delta}{=} P_m^*$. Then we can consider

the family of finite monoids M in which for each idempotent e , $eM_e \cup M_e e = e$. This family is precisely the family of \mathcal{I} -trivial monoids. The generalization of Figure 6 is shown in Figure 7. In general, characterizations of the languages corresponding to the monoids of Figure 7 are not known, except for the \mathcal{I} -trivial case. However, for a two-letter alphabet these languages are generalizations of definite, reverse definite and generalized definite languages, where $\mathcal{F}_\oplus \cup \mathcal{C}_\oplus$ is used instead of $\mathcal{F} \cup \mathcal{C}$.

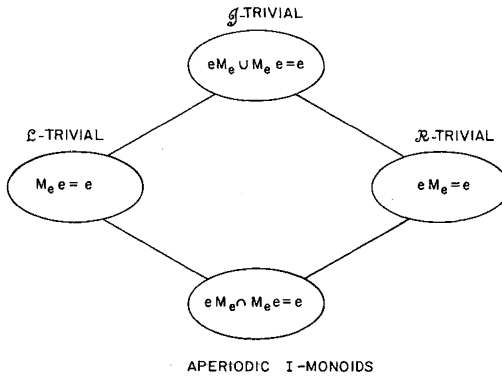


Figure 7.
Aperiodic I-monoids.

The correspondences between the families of languages and monoids or semigroups discussed above are examples of a more general result of Eilenberg [10]. A family of finite monoids is an M -variety if, and only if, it is closed under:

- 1) the operation of taking submonoids;
 - 2) homomorphisms;
- and
- 3) finite direct products.

Let \mathcal{X}_A be the family of all regular languages over alphabet A (subsets of A^*), and let $\mathcal{X} = \{ \mathcal{X}_A; A \text{ is a finite alphabet} \}$. The class \mathcal{X} is called a a^* -variety if, and only if:

- 1) $L \in \mathcal{X}_A$ implies $A^* - L \in \mathcal{X}_A$;
- 2) $L, L' \in \mathcal{X}_A$ implies $L \cap L' \in \mathcal{X}_A$;
- 3) $L \in \mathcal{X}_A$ and $a \in A$ implies

$$\{w; w \in A^*, aw \in L\} \in \mathcal{X}_A \quad \text{and} \quad \{w; w \in A^*, wa \in L\} \in \mathcal{X}_A$$

- 4) If $f: B^* \rightarrow A^*$ is a homomorphism of monoids (where A and B are finite alphabets), and if $L \in \mathcal{X}_A$ then $Lf^{-1} \in \mathcal{X}_B$.

Thus a $*$ -variety is closed under Boolean operations, "removal of a letter" and inverse homomorphisms.

The basic result is that to each $*$ -variety of languages corresponds an M -variety of monoids and vice versa.

In a similar way, if we consider subsets of A^+ instead of A^* , we obtain a $^+$ -variety, and, if we consider finite semigroups instead of monoids, we obtain an S -variety. Again $^+$ -varieties and S -varieties correspond.

11. AUTOMATA

We have already mentioned the following correspondences:

Regular Languages \leftrightarrow Finite Automata
Aperiodic Languages \leftrightarrow Permutation-Free Automata
Definite Languages \leftrightarrow Definite Automata

A characterization of depth-one languages in terms of automata has been found by Simon [23]. We state the result briefly.

A chain-reset is a finite automaton in which the set Q of states can be linearly ordered, say, q_0, q_1, \dots, q_m , in such a way that, for each $q_i \in Q - \{q_m\}$, the next state under any letter of the alphabet is either q_i or q_{i+1} , and for all letters the next state of q_m is q_m . Then $L \in \alpha_{m,k}$, for $m, k \geq 1$, iff the reduced automaton \mathcal{U}_L recognizing L can be covered by a cascade product of two automata \mathcal{U}_1 and \mathcal{U}_2 , where \mathcal{U}_1 is $(k-1)$ -definite and \mathcal{U}_2 is a parallel connection of chain resets with at most $m+1$ states. For more details see [23].

REFERENCES

1. E. BIERMAN, *Realization of Star-Free Events*, M.A.Sc. Thesis, Department of Electrical Engineering, University of Waterloo, Waterloo, Ont., Canada, 1971.
2. J. A. BRZOZOWSKI, *Canonical Regular Expressions and Minimal State Graphs for Definite Events*, Mathematical Theory of Automata, New York, 1962, pp. 529-561, Brooklyn, Polytechnic Institute of Brooklyn, 1963 (Symposia Series, 12).
3. J. A. BRZOZOWSKI, *Run Languages*, Bericht Nr. 87, Institut für Rechner- und Programstrukturen, Gesellschaft für Mathematik und Datenverarbeitung mbH, Bonn, Germany, July 1975, 17 pp.
4. J. A. BRZOZOWSKI, *On aperiodic 1-monoids*, Research Report CS-75-28, Computer Science Department, University of Waterloo, Waterloo, Ont., Canada, November 1975, 18 pp.
5. J. A. BRZOZOWSKI, K. CULIK II, and A. GABRIELIAN, *Classification of Non-counting Events*, J. Computer and System Sc., Vol. 5, 1971, pp. 41-53.
6. J. A. BRZOZOWSKI and I. SIMON, *Characterizations of Locally Testable Events*, Discrete Mathematics, Vol. 4, 1973, pp. 243-271.
7. N. CHOMSKY and M. P. SCHÜTZENBERGER, *The Algebraic Theory of Context-Free Languages*, Computer Programming and Formal Systems, edited by

- P. BRAFFORT and D. HIRSCHBERG, pp. 118-161, Amsterdam, North Holland Publishing Company, 1963.
8. R. S. COHEN and J. A. BRZozowski, *Dot-Depth of Star-Free Events*, J. Computer & System Sc., Vol. 5, 1971, pp. 1-16.
 9. S. EILENBERG, *Automata, Languages, and Machines*, Vol. A, New York, Academic Press, 1974 (Pure and Applied Mathematics Series, 59).
 10. S. EILENBERG, *Automata, Languages and Machines*, Vol. B, New York, Academic Press, 1976.
 11. A. GINZBURG, *About Some Properties of Definite, Reverse Definite and Related Automata*, I.E.E.E. Trans. Electronic Computers EC-15, 1966, pp. 806-810.
 12. S. C. KLEENE, *Representation of Events in Nerve Nets and Finite Automata*, Automata Studies, edited by C.E. SHANNON and J. MCCARTHY, pp. 3-41, Princeton, Princeton University Press, 1954 (Annals of Mathematics Studies, 34).
 13. R. MCNAUGHTON, *Algebraic Decision Procedures for Local Testability*, Math. Systems Theory, Vol. 8, 1974, pp. 60-76.
 14. R. MCNAUGHTON and S. PAPERT, *Counter-Free Automata*, Cambridge, The M.I.T. Press, 1971 (MIT Research Monographs, 65).
 15. YU. T. MEDVEDEV, *On the Class of Events Representable in a Finite Automaton (translated from Russian)*, Sequential Machines-Selected Papers, edited by E.F. MOORE, Reading, Mass., Addison-Wesley, 1964.
 16. A. R. MEYER, *A Note on Star-Free Events*, J. Assoc. Comp. Machin., Vol. 16, 1969, pp. 220-225.
 17. M. PERLES, O. RABIN and E. SHAMIR, *The Theory of Definite Automata*, I.E.E.E. Trans. Electronic Computers EC-12, 1963, pp. 233-143.
 18. D. PERRIN, *Sur certains semigroupes syntaxiques*, Séminaires de l'I.R.I.A. Logiques et Automates, 1971, pp. 169-177.
 19. M. P. SCHÜTZENBERGER, *On Finite Monoids Having Only Trivial Sub-groups*, Inform. and Control, Vol. 8, 1965, pp. 190-194.
 20. M. P. SCHÜTZENBERGER, *On a Family of Sets Related to McNaughton's L-Language*, Automata Theory, edited by E.R. CAIANIELLO, pp. 320-324, New York, Academic Press, 1966.
 21. M. P. SCHÜTZENBERGER, *Sur le produit de concaténation non ambigu*, (to appear in *Semigroup Forum*).
 22. M. STEINBY, *On Definite Automata and Related Systems*, Ann. Acad. Scient. Fennicae, series A.I., 1969, No. 444, 57 pp.
 23. I. SIMON, *Hierarchies of Events With Dot-Depth One*, Ph.D. Thesis, Dept. of Applied Analysis & Computer Science, University of Waterloo, Waterloo, Ont., Canada, 1972.
 24. I. SIMON, *Piecewise Testable Events*, 2nd GI-Professional Conference on Automata Theory and Formal Languages, Kaiserslautern, Germany, May 1975. (To appear in *Lecture Notes in Computer Science*, Springer-Verlag, Berlin).
 25. Y. ZALCSTEIN, *Locally Testable Languages*, J. Computer and System Sc., Vol. 6, 1972, pp. 151-167.
 26. Y. ZALCSTEIN, *Locally Testable Semigroups*, Semigroup Forum, Vol. 5, 1973, pp. 216-227.
 27. Y. ZALCSTEIN, *Syntactic Semigroups of Some Classes of Star-Free Languages*, Automata, Languages and Programming, Proceedings of a Symposium, Rocquencourt, 1972, pp. 135-144, Amsterdam, North-Holland Publishing Company, 1973.