

MÉMOIRES DE LA S. M. F.

WOLFGANG THOMAS

**An application of the Ehrenfeucht-Fraïssé game
in formal language theory**

Mémoires de la S. M. F. 2^e série, tome 16 (1984), p. 11-21

http://www.numdam.org/item?id=MSMF_1984_2_16__11_0

© Mémoires de la S. M. F., 1984, tous droits réservés.

L'accès aux archives de la revue « Mémoires de la S. M. F. » (<http://smf.emath.fr/Publications/Memoires/Presentation.html>) 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/>

AN APPLICATION OF THE EHRENFUCHT-FRAISSÉ GAME
IN FORMAL LANGUAGE THEORY

Wolfgang Thomas

Abstract A version of the Ehrenfeucht-Fraissé game is used to obtain a new proof of a hierarchy result in formal language theory: It is shown that the concatenation hierarchy ("dot-depth hierarchy") of star-free languages is strict.

Résumé Une version du jeu de Ehrenfeucht-Fraissé est appliquée pour obtenir une nouvelle preuve d'un théorème dans la théorie des langages formels: On montre que la hiérarchie de concaténation ("dot-depth hierarchy") des langages sans étoile est stricte.

1. Introduction.

The present paper is concerned with a connection between formal language theory and model theory. We study a hierarchy of formal languages (namely, the dot-depth hierarchy of star-free regular languages) using logical notions such as quantifier complexity of first-order sentences. In this context we apply a form of the Ehrenfeucht-Fraissé game which serves to establish the elementary equivalence between structures with respect to sentences of certain prefix types.

The class of star-free regular languages is of a very basic nature: It consists of all languages (= word-sets) over a given alphabet A which can be obtained from the finite languages by finitely many applications of boolean operations and the concatenation product. (For technical reasons we consider only nonempty words over A , i.e.

languages $L \subset A^+$; in particular, the complement operation is applied w.r.t. A^+ .) General references on the star-free regular languages are McNaughton-Papert (1971), Chapter IX of Eilenberg (1976), or Pin (1984b).

A natural classification of the star-free regular languages is obtained by counting the "levels of concatenation" which are necessary to build up such a language: For a fixed alphabet A , let

$$\begin{aligned} B_0 &= \{L \subset A^+ \mid L \text{ finite or cofinite}\}, \\ B_{k+1} &= \{L \subset A^+ \mid L \text{ is a boolean combination of languages} \\ &\quad \text{of the form } L_1 \cdot \dots \cdot L_n \ (n \geq 1) \text{ with } L_1, \dots, L_n \in B_k\}. \end{aligned}$$

The language classes B_0, B_1, \dots form the so-called dot-depth hierarchy (or: Brzozowski hierarchy), introduced by Cohen/Brzozowski (1971). In the framework of semigroup theory, Brzozowski/Knast (1978) showed that the hierarchy is infinite (i.e. that $B_k \not\supset B_{k-1}$ for $k \geq 1$). The aim of the present paper is to give a new proof of this result, based on a logical characterization of the hierarchy that was obtained in Thomas (1982). The present proof does not rely on semigroup-theory; instead, an intuitively appealing model-theoretic technique is applied: the Ehrenfeucht-Fraissé game.

Let us first state the mentioned characterization result, taking $A = \{a, b\}$. One identifies any word $w \in A^+$, say of length n , with a "word model"

$$w = (\{1, \dots, n\}, <, \min, \max, S, P, Q_a, Q_b)$$

where the domain $\{1, \dots, n\}$ represents the set of positions of letters in the word w , ordered by $<$, where \min and \max are the first and the last position, i.e. $\min = 1$ and $\max = n$, S and P are the successor and predecessor function on $\{1, \dots, n\}$ with the convention that $S(\max) = \max$ and $P(\min) = \min$, and Q_a, Q_b are unary predicates over $\{1, \dots, n\}$ containing the positions with letter a, b respectively. (Sometimes it is convenient to assume that the position-sets of two words u, v are disjoint; then one takes any two nonoverlapping segments of the integers as the position-sets of u and v .) Let L be the first-order language with equality and nonlogical symbols $<, \min, \max, S, P, Q_a, Q_b$. Then the satisfaction of an L -sentence φ in a word w

(written: $w \models \varphi$) can be defined in a natural way, and we say that $L \subset A^+$ is defined by the L -sentence φ if $L = \{w \in A^+ \mid w \models \varphi\}$.

For example, the language $L = (ab)^+$ is defined by

$$Q_a \min \wedge Q_b \max \wedge \forall y (y < \max \rightarrow (Q_a y \leftrightarrow Q_b S(y))) .$$

As usual, a Σ_k -formula is a formula in prenex normal form with a prefix consisting of k alternating blocks of quantifiers, beginning with a block of existential quantifiers. A $B(\Sigma_k)$ -formula is a boolean combination of Σ_k -formulas.

1.1 Theorem. (Thomas (1982)). Let $k > 0$. A language $L \subset A^+$ belongs to B_k iff L is defined by a $B(\Sigma_k)$ -sentence of L .

For the formalization of properties of words the symbols \min, \max, S, P are convenient. But of course they are definable in the restricted first-order language L_0 with the nonlogical constants $<, Q_a, Q_b$ alone. Indeed, we have:

1.2 Lemma. Let $k > 0$. If $L \subset A^+$ is defined by a $B(\Sigma_k)$ -sentence of L , then L is defined by a $B(\Sigma_{k+1})$ -sentence of L_0 .

Proof. The quantifier-free kernel of a Σ_k -formula φ of L can be expressed both by a Σ_2 - and a Π_2 -formula of L_0 . For example, $Q_a S(\min)$ is expressible in the following two ways:

$$(+)\quad \exists y (y = S(\min) \wedge Q_a y), \quad \forall y (y = S(\min) \rightarrow Q_a y)$$

where $y = S(\min)$ is rewritten as a Π_1 -formula of L_0 using

$$\begin{aligned} x = \min &\leftrightarrow \forall z (x = z \vee x < z), \quad x = \max \leftrightarrow \forall z (z = x \vee z < x) \\ S(x) = y &\leftrightarrow (x = \max \wedge x = y) \vee (x < y \wedge \forall z \neg (x < z \wedge z < y)). \end{aligned}$$

Hence we obtain a Σ_{k+1} -sentence of L_0 which is equivalent (in all word-models) to φ by applying one of the two definitions in (+), depending on the case whether the innermost quantifier-block of φ is existential or universal.

We mention without proof that (for $k > 0$) the $B(\Sigma_k)$ -sentences of L_0 define exactly those languages $L \subset A^+$ which occur on the k -th level of another hierarchy of star-free regular languages, introduced by

Straubing (1981). For details concerning the Straubing hierarchy and its relation to the Brzozowski hierarchy cf. Pin (1984a,b). The proof to be given below also shows that the Straubing hierarchy is infinite.

2. The Example Languages

In order to show that $\mathcal{B}_k \supsetneq \mathcal{B}_{k-1}$ for $k \geq 1$, we introduce "example languages" L_k, L_k^+, L_k^- over $A = \{a, b\}$.

Let $|w|_a$ (resp. $|w|_b$) denote the number of occurrences of the letter a (resp. b) in w , and define the weight $\|w\|$ of a word w by

$$\|w\| = |w|_a - |w|_b.$$

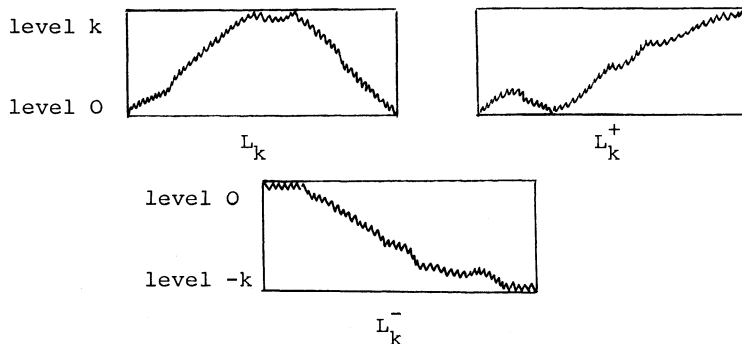
In the sequel we write $v \subseteq w$ if the word v is an initial segment (left factor) of w . Let

$$L_k = \{w \in A^+ \mid \|w\| = 0, \forall v \subseteq w \ 0 \leq \|v\| \leq k, \exists v \subseteq w \ \|v\| = k\},$$

$$L_k^+ = \{w \in A^+ \mid \|w\| = k, \forall v \subseteq w \ 0 \leq \|v\| \leq k\},$$

$$L_k^- = \{w \in A^+ \mid \|w\| = -k, \forall v \subseteq w \ -k \leq \|v\| \leq 0\}.$$

To obtain a more intuitive picture of these languages, it is useful to represent the letter a by the stroke $/$ and b by \backslash . Then the word $abababa$, for example, is represented by $\wedge\wedge\wedge\wedge\wedge$. Thus L_k contains all words whose "graph" has the following properties: It ends on the same level where it starts ("level 0"), it is confined to level 0 and the next k levels, and it assumes the k -th level at least once. Similarly for L_k^+, L_k^- . The "typical shape" of words in L_k, L_k^+, L_k^- is indicated in the following diagrams:



We now state the main result:

2.1 Theorem. For all $k \geq 1$: $L_k \in \mathcal{B}_k - \mathcal{B}_{k-1}$.

The proof is split into lemmas 2.2 and 2.3.

2.2 Lemma. For all $k \geq 1$: $L_k \in \mathcal{B}_k$.

Proof. By induction on k we show that $L_k, L_k^+, L_k^- \in \mathcal{B}_k$. Concerning $k=1$, it is clear that $L_1 = (ab)^+$, $L_1^+ = (ab)^+a$, $L_1^- = b(ab)^+$; hence we can define

$$\begin{aligned} L_1 & \text{ by } (aA^* \cap A^*b) - (A^*aaA^* \cup A^*bbA^*), \\ L_1^+ & \text{ by } (aA^* \cap A^*a) - (A^*aaA^* \cup A^*bbA^*), \\ L_1^- & \text{ by } (bA^* \cap A^*b) - (A^*aaA^* \cup A^*bbA^*). \end{aligned}$$

Observing that, e.g., $A^*aaA^* = aa \cup aaA^+ \cup A^+aa \cup A^+aaA^+$, we see that all three languages belong to \mathcal{B}_1 . - Similarly one obtains, for $k \geq 1$,

$$\begin{aligned} L_{k+1} & = (L_k^+aA^* \cap A^*bL_k^-) - (A^*aL_k^+aA^* \cup A^*bL_k^-bA^*), \\ L_{k+1}^+ & = (L_k^+aA^* \cap A^*aL_k^+) - (A^*aL_k^+aA^* \cup A^*bL_k^-bA^*), \\ L_{k+1}^- & = (L_k^-bA^* \cap A^*bL_k^-) - (A^*aL_k^+aA^* \cup A^*bL_k^-bA^*). \end{aligned}$$

By induction hypothesis, $L_k, L_k^+, L_k^- \in \mathcal{B}_k$; hence, using the elimination of A^* as above, we have $L_{k+1}, L_{k+1}^+, L_{k+1}^- \in \mathcal{B}_{k+1}$.

2.3 Lemma. For all $k \geq 1$: $L_k \notin \mathcal{B}_{k-1}$.

Proof. For $k=1$, the result is clear since $(ab)^+$ is neither finite nor cofinite. By 1.1, it suffices to show for $k \geq 2$ that L_k is not defined by a $B(\Sigma_{k-1})$ -sentence of L . Using 1.2, it is sufficient to prove:

(*) For every $k \geq 2$: L_k is not defined by a $B(\Sigma_k)$ -sentence of L_0 .

Let us write

$u \equiv_n^k v$ iff u and v satisfy the same $B(\Sigma_k)$ -sentences of L_0 in which only prefixes with $\leq n$ quantifiers occur.

We shall verify, for $k \geq 1$, the claim

(*)_k For every $n \geq k$ there are words $u \in L_k, v \notin L_k$ with $u \equiv_n^k v$.

Then in particular for any $k \geq 2$ and $n \geq k$, a $B(\Sigma_k)$ -sentence of L_0 in which only prefixes with $\leq n$ quantifiers occur cannot define L_k , and hence (*) is proved.

The words u, v required in $(*)_k$ for given n will be denoted u_n^k, v_n^k . Together with auxiliary words w_n^k they are defined as follows:

$$u_n^1 = (ab)^{2^n}, \quad v_n^1 = u_n^1 a u_n^1, \quad w_n^1 = u_n^1 b u_n^1,$$

$$u_n^{k+1} = \left(v_n^k w_n^k \right)^{2^n}, \quad v_n^{k+1} = u_n^{k+1} a u_n^{k+1}, \quad w_n^{k+1} = u_n^{k+1} b u_n^{k+1}.$$

(To distinguish superscripts from exponents, the latter are applied only to words in brackets.) The graphs of the first words look as follows (where $n = 2$):

$$u_n^1: \text{ [wavy line] }, v_n^1: \text{ [wavy line] }, w_n^1: \text{ [wavy line]}$$

$$u_n^2: \text{ [wavy line] }$$

From the definition it is immediate that $u_n^k \in L_k, v_n^k \in L_k$. Hence the proof of the main theorem 2.1 is completed when we have shown

$$(**) \quad u_n^k \equiv_n^k v_n^k$$

for $1 \leq k \leq n$. A proof is given in the next section.

3. The Ehrenfeucht-Fraïssé Game G_m .

For the proof that two words are \equiv_n^k -equivalent (as required in (**) above) it is convenient to consider a slight refinement of this notion.

For a sequence $\bar{m} = (m_1, \dots, m_k)$ of positive integers, where $k \geq 0$, let $\text{length}(\bar{m}) = k$ and $\text{sum}(\bar{m}) = m_1 + \dots + m_k$. The set of \bar{m} -formulas (of L_0) is defined by induction on $\text{length}(\bar{m})$: If $\text{length}(\bar{m}) = 0$, it is the set of quantifier-free L_0 -formulas; and for $\bar{m} = (m, m_1, \dots, m_k)$, an \bar{m} -formula is a boolean combination of formulas $\exists x_1 \dots x_m \varphi$ where φ is an (m_1, \dots, m_k) -formula. We write $u \equiv_{\bar{m}} v$ if u and v satisfy the same \bar{m} -sentences. Clearly we have:

3.1 Remark. If $u \equiv_{\bar{m}} v$ for all \bar{m} with $\text{length}(\bar{m}) = k$ and $\text{sum}(\bar{m}) = n$, then $u \equiv_n^k v$.

We now describe the Ehrenfeucht-Fraissé game $G_{\bar{m}}(u, v)$ which is useful for showing $\equiv_{\bar{m}}$ -equivalence. (We restrict ourselves here to the case of word-models for L_0 ; however, all considerations could easily be adapted to arbitrary relational structures.)

The Game $G_{\bar{m}}(u, v)$, where $\bar{m} = (m_1, \dots, m_k)$, is played between two players I and II on the word-models u and v ; we assume that the position-sets of u and v are disjoint. We write $<^u$ to denote the $<$ -relation in u ; $Q_a^u, Q_b^u, <^v, Q_a^v, Q_b^v$ are used similarly. A play of the game consists of k moves. In the i -th move player I chooses, in u or in v , a sequence of m_i positions; then player II chooses, in the remaining word (v or u), also a sequence of m_i positions. Before each move, player I has to decide whether to choose his next elements from u or from v . After k moves, by concatenating the position-sequences chosen from u and chosen from v , two sequences $\bar{p} = p_1 \dots p_n$ from u and $\bar{q} = q_1 \dots q_n$ from v have been formed where $n = m_1 + \dots + m_k$. Player II has won the play if the map $p_i \mapsto q_i$ respects $<$ and the predicates Q_a, Q_b (i.e. $p_i <^u p_j$ iff $q_i <^v q_j$, $Q_a^u p_i$ iff $Q_a^v q_i$, $Q_b^u p_i$ iff $Q_b^v q_i$ for $1 \leq i, j \leq n$). Equivalently, the two subwords in u and v given by the position-sequences \bar{p} and \bar{q} should coincide. If there is a winning strategy for II in the game (to win each play) we say that II wins $G_{\bar{m}}(u, v)$ and write $u \sim_{\bar{m}} v$.

The standard Ehrenfeucht-Fraissé game is the special case of $G_{\bar{m}}(u, v)$ where $\bar{m} = (1, \dots, 1)$. (For a detailed discussion cf. Rosenstein (1982).) If $\text{length}(\bar{m}) = k$ and $\bar{m} = (1, \dots, 1)$ we write $G_k(u, v)$ instead of $G_{\bar{m}}(u, v)$ and $u \sim_k v$ instead of $u \sim_{\bar{m}} v$. Note that in this case the \bar{m} -formulas are (up to equivalence) just the formulas of quantifier-depth k . In the familiar form the Ehrenfeucht-Fraissé Theorem states (for the case of word-models) that u and v satisfy the same L_0 -sentences of quantifier-depth k iff $u \sim_k v$. An analogous proof yields the result for \bar{m} -sentences and $\sim_{\bar{m}}$ (cf. Fraissé (1972), where the terminology of partial isomorphisms is used instead of game-theoretical notions):

3.2 Theorem. For all $\bar{m} = (m_1, \dots, m_k)$ with $k > 0$ and $m_i > 0$ for $i = 1, \dots, k$, we have $u \equiv_{\bar{m}} v$ iff $u \sim_{\bar{m}} v$.

Hence, in view of 3.1, we can prove the claim (**) of the preceding section (and thus the main result 2.1) by showing

3.3 Lemma. For $0 < k \leq n$ and any \bar{m} with length $(\bar{m}) = k$ and $\text{sum}(\bar{m}) = n$,
 $u_n^k \sim_{\bar{m}}^k v_n^k$ and $u_n^k \sim_{\bar{m}}^k w_n^k$.

As a preparation for the proof we state some basic properties of $\sim_{\bar{m}}$ and \sim_n :

3.4 Lemma.

- (a) $\sim_{\bar{m}}$ is an equivalence relation.
- (b) If $n \geq \text{sum}(\bar{m})$ and $u \sim_n v$, then $u \sim_{\bar{m}} v$.
- (c) If $u \sim_{\bar{m}} v$ and $u' \sim_{\bar{m}} v'$, then $uu' \sim_{\bar{m}} vv'$.

Parts (a), (b) are immediate from the definition of $G_n(u, v)$ and $G_{\bar{m}}(u, v)$. For the proof of (c) note that player II can combine the two given winning strategies on u, v and on u', v' in the obvious manner to obtain a winning strategy on uu', vv' : As far as the initial segments u and v are concerned, the first given strategy is to be used, similarly for the final segments u', v' the second given strategy.

The following lemma is a familiar exercise on the game:

3.5 Lemma. If $m, m' \geq 2^n - 1$, then $(w)^m \sim_n (w)^{m'}$.

Proof. Consider the natural decomposition of $u = (w)^m$ and $v = (w)^{m'}$ into w -segments. Before each move we have in u and v certain w -segments in which positions have been chosen, and others where no positions have been chosen. Call a maximal segment of succeeding w -segments without chosen positions a gap. (A gap may be empty.) Before each move there is a natural correspondence between the gaps in u and v (given by their order). II should play according to what we call the 2^i -strategy, namely guarantee the following condition before each move: When i elements are still to be chosen by both players, two corresponding gaps should both consist of any number $\geq 2^i - 1$ of w -segments, or else should both consist of the same number ($< 2^i - 1$) of w -segments. By induction on $n-i$ it is easy to see that II always can choose his w -segment in this manner (cf. Rosenstein (1982), p. 99); of course, inside his w -segment, II should pick exactly that position which matches the position chosen by I in the corresponding w -segment.

Since any word u_n^k as defined in §2 is of the form $(w)^{2^n}$, we note as a consequence of 3.5:

3.6 Remark. For $1 \leq k \leq n$: $u_n^k \sim_n u_n^k u_n^k$.

We now turn to the

Proof of 3.3. By induction on k we show $u_n^k \sim_{\bar{m}} v_n^k$ and $u_n^k \sim_{\bar{m}} w_n^k$ for any \bar{m} with $\text{length}(\bar{m}) = k$ and $\text{sum}(\bar{m}) \leq n$.

If $k = 1$ we deal with the game involving one move in which $\leq n$ elements are chosen by both players. Let us consider

$$u_n^1 = (ab)^{2^n}, \quad v_n^1 = (ab)^{2^n} a (ab)^{2^n}.$$

Since in both words u_n^1 and v_n^1 all possible words over $\{a, b\}$ of length n occur as subwords, any subword specified by I through his choice of n positions in one word can also be realized by II in the remaining word by n corresponding positions. Hence there is a winning strategy for II. The proof for u_n^1 and w_n^1 is analogous.

In the induction step we write u for u_n^k and consider the words

$$u_n^{k+1} = (uauubu)^{2^n}, \quad v_n^{k+1} = (uauubu)^{2^n} a (uauubu)^{2^n}.$$

Given a sequence (m, \bar{m}) with $\text{length}(m, \bar{m}) = k + 1$ and $\text{sum}(m, \bar{m}) \leq n$, we have to show $u_n^{k+1} \sim_{(m, \bar{m})} v_n^{k+1}$, using as induction hypothesis

$$(a) \quad u \sim_{\bar{m}} uau \quad (= v_n^k), \quad (b) \quad u \sim_{\bar{m}} ubu \quad (= w_n^k).$$

(In an analogous manner it will be possible to show $u_n^{k+1} \sim_{(m, \bar{m})} w_n^{k+1}$.)

In order to verify $u_n^{k+1} \sim_{(m, \bar{m})} v_n^{k+1}$, it is convenient to apply 3.4(a), (b) and consider two different words instead which are \sim_n -equivalent to u_n^{k+1} , v_n^{k+1} respectively: Instead of u_n^{k+1} we take

$$(1) \quad (uauubu)^{2^n} uauubu (uauubu)^{2^n}$$

which is \sim_n -equivalent to u_n^{k+1} by 3.5. Concerning v_n^{k+1} , we use 3.6 in order to duplicate (several times) the u -segments next to the central letter a there; thus we obtain the \sim_n -equivalent word

$$(2) \quad (uauubu)^{2^n} uau (u)^{m+1} (uauubu)^{2^n}.$$

For the proof of $(1) \sim_{(m, \bar{m})} (2)$ we distinguish the two cases that I first picks m positions from (1) or I first picks m positions from (2).