Mémoires de la S. M. F.

WOLFGANG THOMAS

An application of the Ehrenfeucht-Fraisse 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.

\mathcal{N} umdam

Article numérisé dans le cadre du programme Numérisation de documents anciens mathématiques http://www.numdam.org/ Société Mathématique de France 2° série, mémoire n° 16, 1984, p. 11-21

AN APPLICATION OF THE EHRENFEUCHT-FRAISSE GAME IN FORMAL LANGUAGE THEORY

Wolfgang Thomas

<u>Abstract</u> A version of the Ehrenfreucht-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.

<u>Résumé</u> 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.

0037-9484/84 03 11 11/\$ 3.10/ © Gauthier-Villars

W. THOMAS

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).

Į

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

$$B_{O} = \{ L \subset A^{+} | L \text{ finite or cofinite} \},\$$

$$B_{k+1} = \{ L \subset A^{+} | L \text{ is a boolean combination of languages} \\ \text{ of the form } L_{1} \cdot \ldots \cdot L_{n} (n \ge 1) \text{ with } L_{1} \cdot \ldots \cdot L_{n} \in B_{k} \}.$$

The language classes B_0, B_1, \ldots 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 $B_k \stackrel{?}{\downarrow} B_{k-1}$ for $k \ge 1$). 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"

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

where the domain $\{1, \ldots, n\}$ 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 $\{1, \ldots, n\}$ with the convention that S(max) = max and P(min) = min, and Q_a, Q_b are unary predicates over $\{1, \ldots, n\}$ 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

EHRENFEUCHT-FRAISSE GAME

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

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

I

 $Q_{a}\min \wedge Q_{b}\max \wedge \forall y (y < \max \rightarrow (Q_{a}y \leftrightarrow Q_{b}S(y))) .$

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(Σ_k)-formula is a boolean combination of Σ_k -formulas.

<u>1.1</u> 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:

<u>1.2</u> 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 .

<u>Proof.</u> 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:

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

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

 $\begin{array}{l} x = \min \iff \forall z \ (x = z \ \lor \ x < z), \ x = \max \iff \forall z \ (z = x \ \lor \ z < x) \\ S \ (x) = y \iff (x = \max \land \ x = y) \ \lor \ (x < y \land \forall z \neg (x < z \land \ z < y)). \end{array}$

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

W. THOMAS

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 $B_k \stackrel{\neg}{=} B_{k-1}$ for $k \ge 1$, 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 $||w|| \circ f a$ word w by

$$\|\mathbf{w}\| = \|\mathbf{w}\|_{a} - \|\mathbf{w}\|_{b}$$

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

$$\begin{split} \mathbf{L}_{\mathbf{k}} &= \{ \mathbf{w} \in \mathbf{A}^{+} | \| \mathbf{w} \| = 0, \ \forall \mathbf{v} \underline{-} \mathbf{w} \ \mathbf{0} \leq \| \mathbf{v} \| \leq \mathbf{k}, \ \exists \mathbf{v} \underline{-} \mathbf{w} \| \mathbf{v} \| = \mathbf{k} \} , \\ \mathbf{L}_{\mathbf{k}}^{+} &= \{ \mathbf{w} \in \mathbf{A}^{+} | \| \mathbf{w} \| = \mathbf{k}, \ \forall \mathbf{v} \underline{-} \mathbf{w} \ \mathbf{0} \leq \| \mathbf{v} \| \leq \mathbf{k} \} , \\ \mathbf{L}_{\mathbf{k}}^{-} &= \{ \mathbf{w} \in \mathbf{A}^{+} | \| \mathbf{w} \| = -\mathbf{k}, \ \forall \mathbf{v} \underline{-} \mathbf{w} - \mathbf{k} \leq \| \mathbf{v} \| \leq \mathbf{0} \} . \end{split}$$

To obtain a more intuitive picture of these languages, it is useful to represent the letter a by the stroke \checkmark and b by \searrow . Then the word abababa, for example, is represented by \checkmark . 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:

