# On the regularity of languages on a binary alphabet generated by copying systems

## Daniel P. Bovet

*Dipartimento di Scienze dell'Informazione, Università di Roma "La Sapienza", 00185 Roma, Italy*

## Stefano Varricchio

*Dipartimento de Matematica, Università dell'Aquila, 67100 L'Aquila, Italy*

*Abstract*

Bovet, D.P. and S. Varricchio, On the regularity of languages on a binary alphabet generated by copying systems, Information Processing Letters 44 (1992) 119–123.

Let $A$ be a binary alphabet and $\vdash_\pi^*$ the derivation relation associated to the semi-Thue system $\pi = \{(x, xx) \mid x \in A^*\}$. We prove that $\vdash_\pi^*$ is a well quasi order in $A^*$. As a consequence of this we show that the languages which are upwards closed with respect to $\vdash_\pi^*$ are all regular.

*Keywords*: Formal languages; copying systems; combinatorics on words; well quasi orders

## 1. Introduction

Copying systems and languages generated by them were introduced in [4] by Ehrenfeucht and Rozenberg. In their papers it is proved that when the alphabet has cardinality at least three then such languages are not, in general, regular. In this paper, following an approach started in [2] we prove that in the case of a binary alphabet, the languages generated by copying systems are all regular.

At the base of our construction there is the notion of well quasi order (w.q.o.). There exist different equivalent definitions of this concept

which was often rediscovered by many authors and there is a large literature on this subject. We recall the division ordering which is a well quasi order, as it was proved by Higman in his famous paper [5]. The relevance of well quasi orders in Computer Science and Automata Theory is due to the following theorem of Ehrenfeucht et al. [3]: *A language L of a finitely generated free monoid is regular if and only if it is upwards closed with respect to a monotone well quasi order*. This result has been used in [2] in order to give many interesting regularity conditions for languages.

An important class of quasi orders is obtained considering the derivation relations associated to rewriting systems. Thus an interesting problem has become that of finding rewriting systems whose associated derivation relation is a w.q.o. In this paper we prove that in the case of a binary

*Correspondence to*: Professor S. Varricchio, Dipartimento di Matematica, Università dell'Aquila, 67100 L'Aquila, Italy.

alphabet the copying relation introduced in [4], is a monotone w.q.o., and thus all languages which are upwards closed with respect to this relation are regular.

In the sequel $A$ will denote a finite set or alphabet, and $A^+$ (resp. $A^*$) the *free semigroup* (resp. *free monoid*) over $A$. The elements of $A$ are called *letters* and those of $A^*$ *words*. The identity element of $A^*$ is denoted by $\Lambda$. For any word $w$, $|w|$ denotes its *length*. A word $u$ is a *factor* of the word $w$ (resp. *prefix, suffix*) if $w \in A^*uA^*$ (resp. $w \in uA^*$, $w \in A^*u$). For any $w \in A^*$, $F(w)$ denotes the set of its factors. A language $L$ of $A^*$ is called *regular* if $L$ is recognized by a finite state automaton.

A binary relation $\leqslant$ in a set $T$ is a *quasi order* if it is reflexive and transitive. If for any $a,b \in T$, $a \leqslant b$ and $b \leqslant a$ implies $a = b$, then $\leqslant$ is called *a partial order*. Given a quasi order $\leqslant$, we set $a < b$ if $a \leqslant b$ and $a \neq b$. An *anti-chain* is a set $C \subseteq T$ such that for any $a,b \in C$, $a \neq b$, one has $a \nleqslant b$. An infinite sequence $\{x_n\}_{n > 0}$ of elements of $T$ is a *bad sequence* if for any $i,j$, $i < j$, one has $x_i \nleqslant x_j$. An infinite antichain $C$ contains trivially a bad sequence, but in general a bad sequence $\{x_n\}_{n > 0}$ does not contain an infinite anti-chain, since for some $i,j$, $i < j$, one could have $x_j \leqslant x_i$. We say that a quasi order $\leqslant$ is a *well quasi order* (w.q.o.) if the two following conditions are verified:

(1) there is no infinite strictly descending sequence (i.e. $\leqslant$ is well founded),

(2) there are not infinitely many elements which are mutually incomparable with respect to $\leqslant$ (i.e. there is no infinite anti-chain).

A w.q.o. in the free monoid $A^*$ is called *monotone* if for any $x,y,u,v \in A^*$, $u \leqslant v$ implies $xuy \leqslant xvz$. We say that a language $L \subseteq A^*$ is *upwards-closed* with respect to $\leqslant$ if for any $x,y \in A^*$, $x \in L$ and $x \leqslant y$ implies $y \in L$.

We recall [1] that a *rewriting system* or *semi-Thue system* is a pair $(A, \pi)$ where $\pi$ is a binary relation on $A^*$, i.e. $\pi \subseteq A^* \times A^*$. Any pair $(x, y) \in \pi$ is called a *production* and denoted by $x \to y$. Let us denote by $\vdash_\pi$ the regular closure of $\pi$, i.e. for any $u,v \in A^*$, $u \vdash_\pi v$ if and only if $\exists (x, y) \in \pi$ and $\exists h,k \in A^*$ such that $u = hxk$ and $v = hyk$. The *derivation relation* $\vdash_\pi^*$ is the

reflexive and transitive closure of $\vdash_\pi$. It is easy to verify that the derivation relation $\vdash_\pi^*$ is a monotone quasi order.

## 2. The main result

In this section $A$ will denote a binary alphabet $A = \{a, b\}$. Let us also denote with $\pi$ the rewriting system $\pi = \{(x, xx) \mid x \in A^*\}$. The derivation relation $\vdash_\pi^*$ has been considered in [4] where it is called *copying relation*. In this section we shall consider a restricted copying relation denoted as $\vdash_{\pi'}^*$ where $\pi' = \{(a, aa),\ (b, bb),\ (ab, abab),\ (ba, baba)\}$. Observe that $\vdash_{\pi'}^* \subseteq \vdash_\pi^*$, thus if we succeed in proving that $\vdash_{\pi'}^*$ is a w.q.o., then also $\vdash_\pi^*$ turns out to be a w.q.o. In fact, if an infinite antichain $C$ exists with respect to $\vdash_\pi^*$, then $C$ would also be an antichain with respect to $\vdash_{\pi'}^*$. The reason why we consider such a restricted rewriting system is that $\pi'$ is the smallest set of rules among those equivalent to $\pi$. We begin with some preliminary lemmas.

**Lemma 2.1.** *Let $\lambda,\mu \in A^*$ such that $aba\lambda \vdash_{\pi'}^* aba\mu$. Then, $ba\lambda \vdash_{\pi'}^* ba\mu$.*

**Proof.** By induction on the length of the derivation.

Consider first a derivation of length 1, that is, $aba\lambda \vdash_{\pi'} aba\mu$. Two subcases may occur depending on whether the derivation makes use of the leftmost symbol $a$ or not. In the second case $(ba\lambda \vdash_{\pi'} ba\mu)$, the assertion follows trivially; in the first case the derivation must be of the following type:

$$(ab)a\lambda \vdash_{\pi'} (abab)a\lambda = aba\mu$$

with $\mu = ba\lambda$. On the other hand, $(ba)\lambda \vdash_{\pi'} (baba)\lambda = ga\mu$. The base of the induction is thus proved.

Let us consider derivations of length $n \sim 1$ of the form

$$aba\lambda \vdash_{\pi'} v_1 \vdash_{\pi'} v_2 \vdash_{\pi'} \cdots \vdash_{\pi'} v_{n-1} \vdash_{\pi'} aba\mu.$$

$$(1)$$

We distinguish the following cases:

(1) The first derivation is obtained by applying $\vdash_{\pi'}$ to the word $a\lambda$. In this case, we can write

$$aba\lambda \vdash_{\pi'} aba\lambda' = v_1,$$

where $a\lambda \vdash_{\pi'} a\lambda'$ and $aba\lambda' \vdash_{\pi'}^* aba\mu$. Since the latter derivation has length $n - 1$, by the induction hypothesis, $ba\lambda' \vdash_{\pi'}^* ba\mu$. On the other hand, $a\lambda \vdash_{\pi'} a\lambda'$ implies $ba\lambda \vdash_{\pi'} ba\lambda'$, thus $ba\lambda \vdash_{\pi'} ba\mu$.

(2) The first derivation is obtained by applying $\vdash_{\pi'}$ to the word $aba$. Let us consider the following subcases:

(a) The first derivation is $a(ba)\lambda \vdash_{\pi'} a(baba)\lambda = v_1$. In such a case, $v_1$ yields $aba\mu$ in $n - 1$ steps and, according to the inductive hypothesis, one has $baba\lambda \vdash_{\mu'}^* ba\mu$. Moreover, $ba\lambda \vdash_{\pi'} baba\lambda$, thus $ba\lambda \vdash_{\pi'}^* ba\mu$.

(b) The first derivation is $(ab)a\lambda \vdash_{\pi'} (abab)a\lambda = v_1 = a(baba)\lambda$. This subcase is equivalent to subcase (a).

(c) The first derivation is $(a)ba\lambda \vdash_{\pi'} (aa)ba\lambda = v_1$. This case cannot occur since $v_1$ cannot yield $aba\mu$ (a word starting with $aa$ can only yield words starting with $aa$).

(d) The first derivation is $ab(a)\lambda \vdash_{\pi'} ab(aa)\lambda = v_1$. This subcase is included in case 1.

(e) The first derivation is $a(b)a\lambda \vdash_{\pi'} a(bb)a\lambda = v_1$. We observe that in the derivation $v_1 = abba\lambda \vdash_{\pi'}^* aba\mu$, there is at least a direct derivation of the first two symbols $ab$ into $abab$. Indeed, if this is not the case, it is easy to see that we could not remove the prefix $abb$ in order to obtain $aba\mu$. Then the derivation (1) may be rewritten as

$$aba\lambda \vdash_{\pi'} abba\lambda \vdash_{\pi'}^* abb\lambda' \vdash_{\pi'} ababb\lambda' \vdash_{\pi'}^* aba\mu.$$

The previous derivation may be reordered in the following one

$$aba\lambda \vdash_{\pi'} ababa\lambda \vdash_{\pi'} ababba\lambda$$
$$\vdash_{\pi'}^* ababb\lambda' \vdash_{\pi'}^* aba\mu.$$

We observe that this derivation has length $n$ and may be dealt with as in subcase (b).    □

**Corollary 2.2.** *Let* $\lambda,\mu \in A^*$ *such that* $aba\lambda \vdash_{\pi'}^*$ *$aba\mu$. Then, for all $k > 0$, $ab^k a\lambda \vdash_{\pi'}^* ab^k a\mu$.*

**Proof.** According to Lemma 2.1, $aba\lambda \vdash_{\pi'}^* aba\mu$ implies $ba\lambda \vdash_{\pi'}^* ba\mu$. After performing $k - 1$ left compositions of $b$'s, we get $b^k a\lambda \vdash_{\pi'}^* b^k a\mu$; the result follows by applying a final left composition of $a$.    □

**Corollary 2.3.** *Let* $\lambda,\mu \in A^*$ *such that* $abab\lambda \vdash_{\pi'}^*$ *$abab\mu$. Then, for all $k > 0$, $aba^k b\lambda \vdash_{\pi'}^* aba^k b\mu$.*

**Proof.** According to Lemma 2.1, $aba(b\lambda) \vdash_{\pi'}^*$ $aba(b\mu)$ implies $ba(b\lambda) \vdash_{\pi'}^* ba(b\mu)$. By applying a second time Lemma 2.1 (and exchanging the role of $a$'s and $b$'s), $ab\lambda \vdash_{\pi'}^* ab\mu$ follows. A series of $(k - 1)$ left compositions with $a$, one with $b$ and one with $a$ yields the result.    □

Clearly for symmetry reasons Lemma 2.1, Corollary 2.2 and Corollary 2.3 hold when $\lambda$ and $\mu$ appear as leftmost word.

**Theorem 2.4.** *The derivation relation* $\vdash_{\pi'}^*$ *is a well quasi order.*

**Proof.** It is easy to verify that $\vdash_{\pi'}^*$ is well founded. Suppose by contradiction that a bad sequence $C = x_1, x_2, \ldots, x_n, \ldots$ exists. Following a classical argument (cf. [2]), we may assume, without loss of generality, that it is *minimal* in the following sense: no bad sequence $y_1, y_2, \ldots, y_n, \ldots$ exists such that for some $j > 0$ one has $x_i = y_i$ for $i = 1, \ldots, j - 1$ and $|y_j| < |x_j|$.

Clearly $C$ includes infinitely many words starting with the same letter, say $a$. We distinguish the following cases:

(1) $C$ includes infinitely many words starting with the same two letters $aa$. Denote with $aay_1, aay_2, \ldots, aay_n, \ldots$ such an infinite subsequence and let $j$ be the integer such that $aay_1 = x_j$. Consider the following infinite sequence $D = z_1, z_2, \ldots, z_n, \ldots$ where $z_i = x_i$ for $i = 1, 2, \ldots,$ $j - 1$ and $z_i = ay_{i-j+1}$ for $i \geqslant j$.

It is easy to verify that $D$ is a bad sequence: if $s$ and $t$ are such that $s < t$ and $s, t \in \{1, \ldots, j - 1\}$, then $z_s \nvdash_{\pi'}^* z_t$ by assumption since $z_s = x_s$ and $z_t = x_t$. If $s \in \{1, \ldots, j - 1\}$ and $t \geqslant j$, then $z_s = x_s$ and $z_t = ay_k$ with $k = t - j + 1$. From $z_s \vdash_{\pi'}^* z_t$, it follows $x_s \vdash_{\pi'}^* ay_k \vdash_{\pi'}^* aay_k$. This contradicts the assumption that $C$ is a bad sequence. Finally if

$s, t > j$, then $z_s = ay_h$ and $z_t = ay_k$ with $h < k$. If $z_s \vdash_{\pi'}^* z_t$, due to left regularity, $aay_h \vdash_{\pi}^* aay_k$ and this contradicts the assumption that $C$ is a bad sequence.

(2) $C$ includes infinitely many words starting with $ab$. If $C$ includes an infinite number of words belonging to $ab^*$, then it also includes infinitely many words terminating with $bb$. This case is symmetric to the previous one. Let us then suppose that $C$ includes an infinite subsequence $ab^{k_1}ay_1, ab^{k_2}ay_2, \ldots, ab^{k_n}ay_n, \ldots$ with $k_i > 1$. One has to consider the following subcases:

(a) $C$ includes an infinite subsequence of elements of the form $ab^{k_i}ay_i$ with $k_i \geqslant 2$ for all $i$. We can always suppose that the $\{k_i\}_{i>1}$ is a nondecreasing sequence.

Let $j$ be the integer such that $ab^{k_1}ay_1 = x_j$. Consider the following infinite sequence $D = z_1, z_2, \ldots, z_n, \ldots$ where $z_i = x_i$ for $i = 1, 2, \ldots, j-1$ and $z_i = abay_{i-j+1}$ for $i \geqslant j$.

It is easy to verify that $D$ is a bad sequence: if $s$ and $t$ are such that $s < t$ and $s, t \in \{1, \ldots, j-1\}$, then $z_s \not\vdash_{\pi'}^* z_t$ by assumption since $z_s = x_s$ and $z_t = x_t$. If $s \in \{1, \ldots, j-1\}$ and $t \geqslant j$, then $z_s = x_s$ and $z_t = abay_p$ with $p = t - j + 1$. From $z_s \vdash_{\pi'}^* z_t$, it follows $x_s \vdash_{\pi'}^* abay_p \vdash_{\pi'}^* ab^{k_p}ay_p$. This contradicts the assumption that $C$ is a bad sequence. Finally if $s, t > j$, then $z_s = abay_p$ and $z_t = abay_q$ with $p < q$. If $z_s \vdash_{\pi'}^* z_t$, due to Corollary 2.2, $ab^{k_p}ay_p \vdash_{\pi}^* ab^{k_p}ay_q$; since $k_p \leqslant k_q$ and $ab^{k_p}ay_q \vdash^* ab^{k_q}ay_q$, then $ab^{k_p}ay_p \vdash_{\pi'}^* ab^{k_q}ay_q$, contradicting the assumption.

(b) $C$ includes an infinite subsequence of elements of the form $abay_i$. Now we have to consider the following subcases:

(i) $C$ includes an infinite sequence of elements of the kind $ababy_i$. Denote with $ababy_1$, $ababy_2, \ldots, ababy_n, \ldots$ such an infinite subsequence and let $j$ be the integer such that $ababy_1 = x_j$. Consider the following infinite sequence $D = z_1, z_2, \ldots, z_n, \ldots$ where $z_i = x_i$ for $i = 1, 2, \ldots, j-1$ and $z_i = aby_{i-j+1}$ for $i \geqslant j$.

It is easy to verify that $D$ is a bad sequence: if $s$ and $t$ are such that $s < t$ and $s, t \in \{1, \ldots, j-1\}$, then $z_x \not\vdash_{\pi'}^* z_t$ by assumption since $z_s = x_s$ and $z_t = x_t$. If $s \in \{1, \ldots, j-1\}$ and $t \geqslant j$, then $z_s = x_s$ and $z_t = aby_k$ with $k = t - j + 1$. From $z_s \vdash_{\pi'}^* z_t$, it follows $x_s \vdash_{\pi'}^* aby_k \vdash_{\pi'} ababy_k$. This contra-

dicts the assumption that $C$ is a bad sequence. Finally if $s, t > j$, then $z_s = aby_h$ and $z_t = aby_k$ with $h < k$. If $z_s \vdash_{\pi'}^* z_t$, due to left regularity, $ababy_h \vdash_{\pi}^* ababy_h$ and this contradicts the assumption that $C$ is a bad sequence.

(ii) $C$ includes an infinite sequence of elements of the kind $abaay_i$. If $C$ contains infinitely many elements $\in abaa_*$, then it also contains infinite elements which terminate with $aa$ and this case is symmetric to case 1. Thus we can suppose that $C$ includes an infinite subsequence of elements of the form $aba^{k_i}by_i$ with $k_i \geqslant 2$ for all $i$. We can always suppose that the $\{k_i\}_{i \geqslant 1}$ is a nondecreasing sequence.

Let $j$ be the integer such that $aba^{k_1}by_1 = x_j$. Consider the following infinite sequence $D = z_1, z_2, \ldots, z_n, \ldots$ where $z_i = x_i$ for $i = 1, 2, \ldots, j-1$ and $z_i = ababy_{i-j+1}$ for $i \geqslant j$.

It is easy to verify that $D$ is a bad sequence: if $s$ and $t$ are such that $s < t$ and $s, t \in \{1, \ldots, j-1\}$, then $z_s \not\vdash_{\pi'}^* z_t$ by assumption since $z_s = x_s$ and $z_t = x_t$. If $s \in \{1, \ldots, j-1\}$ and $t \geqslant j$, then $z_s = x_s$ and $z_t = ababy_p$ with $p = t - j + 1$. From $z_s \vdash_{\pi'}^* z_t$, one derives $x_s \vdash_{\pi'}^* ababy_p \vdash_{\pi'}^* aba^{k_p}by_p$. This contradicts the assumption that $C$ is a bad sequence. Finally if $s, t > j$, then $z_s = ababy_p$ and $z_t = ababy_q$ with $p < q$. If $z_s \vdash_{\pi'}^* z_t$, due to Corollary 2.3, $aba^{k_p}by_p \vdash_{\pi}^* aba^{k_p}by_q$; since $k_p \leqslant k_q$ and $aba^{k_p}by_q \vdash_{\pi}^* aba^{k_q}by_q$, then $aba^{k_p}by_p \vdash_{\pi'}^* aba^{k_q}by_q$ contradicting the assumption.

In all the cases we have considered one derives the existence of a bad sequence $D$ which contradicts the minimality of $C$.  □

**Corollary 2.5.** *The derivation relation $\vdash_{\pi}^*$ is a well quasi order.*

**Proof.** It is evident that $\vdash_{\pi}^*$ is well founded. Moreover, suppose by contradiction that there is an infinite antichain $C$ in $A^*$ with respect to $\vdash_{\pi}^*$. Since $\vdash_{\pi'}^* \subseteq \vdash_{\pi}^*$, $C$ should also be an antichain with respect to $\vdash_{\pi}^*$ and this is a contradiction.  □

We recall now the following important theorem [3] which gives a noteworthy generalization of the Myhill–Nerode Theorem for regular languages.

**Theorem 2.6.** *A language L of a finitely generated free monoid $B^*$ is regular if and only if it is closed with respect to a monotone well quasi order.*

**Corollary 2.7.** *Let $L \subseteq A^*$ be a language which is upwards closed with respect to $\vdash_\pi^*$. Then L is a regular language.*

**Proof.** The statement is a consequence of Corollary 2.5 and Theorem 2.6. $\square$

Let us now consider a free monoid $B^*$ and the copying relation $\vdash_\pi^*$ in $B^*$. For any $w \in B^*$ we consider the set $L_{w,\pi}$ defined by

$$L_{w,\pi} = \{ u \in B^* \mid w \vdash_\pi^* u \}.$$

It has been proved in [2] that, if $w$ is a word containing at least three letters, then $L_{w,\pi}$ is not a regular language. We can conclude by observing that $L_{w,\pi}$ is upwards closed with respect to $\vdash_\pi^*$; thus taking in account of Corollary 2.7 we easily derive the following:

**Corollary 2.8.** *Let $B$ a finite alphabet and $w \in B^*$. Then $L_{w,\pi}$ is regular if and only if $w$ contains at most two letters.*

## References

[1] J. Berstel, *Transductions and Context-Free Languages* (Teubner, Stuttgart, 1979).

[2] A. de Luca and S. Varricchio, Some regularity conditions based on well quasi-orders, in: *Proc. Conf. Latin 92*, Lectures Notes in Computer Science (Springer, Berlin).

[3] A. Ehrenfeucht, D. Haussler and G. Rozenberg, On regularity of context-free languages, *Theoret. Comput. Sci.* **28** (1983) 311–332.

[4] A. Ehrenfeucht and G. Rozenberg, On regularity of languages generated by copying systems, *Discrete Appl. Math.* **8** (1984) 313–317.

[5] G.H. Higman, Ordering by divisibility in abstract algebras, *Proc. London Math. Soc.* **3** (1952) 326–336.

[6] M. Lothaire, *Combinatorics on Words* (Addison-Wesley, Reading, MA, 1983).