



# Shuffle factorization is unique

Jean Berstel<sup>a</sup>, Luc Boasson<sup>b,\*</sup><sup>a</sup>*Institut Gaspard Monge (IGM), Université Marne-la-Vallée, 5, boulevard Descartes, 77454 Marne-la-Vallée Cédex 2, France*<sup>b</sup>*Laboratoire d'informatique algorithmique: fondements et applications (LIAFA), Université Denis-Diderot, 2, place Jussieu, 75251 Paris Cédex 05, France*

## Abstract

We prove that, given a finite set of words  $S$ , there exists at most one (normalized) multiset  $\mathcal{P}$  such that  $S$  is the shuffle of the words in  $\mathcal{P}$ . The multiset  $\mathcal{P}$  is effectively computable.  
© 2002 Elsevier Science B.V. All rights reserved.

## 1. Introduction

Let  $A$  be an alphabet. The *shuffle*  $x \sqcup\sqcup y$  of words  $x, y$  over  $A$  is a binary commutative and associative operation. The paper addresses the following problem: “Given a finite set  $S$  of words over  $A$ , do there exist words  $x_1, x_2, \dots, x_n$  such that

$$S = x_1 \sqcup\sqcup x_2 \sqcup\sqcup \dots \sqcup\sqcup x_n \quad (1)$$

and, in the positive case, is the solution unique?”

A set  $S$  such that Eq. (1) holds, is called a *shuffle set*. Observe that if Eq. (1) holds, then all words in  $S$  have same length and even have same Parikh vector. This means that words in a solution  $x_1, x_2, \dots, x_n$  have bounded length, and thus an exhaustive search answers the question of existence, with a considerable cost. We give here a different method. This method could presumably lead to an efficient algorithm. This will be a topic of later work.

Next, if a solution  $x_1, x_2, \dots, x_n$  exists, there are easy cases where it is not unique, even up to permutation. Consider a single letter  $a$ , and any power  $a^n$  with  $n \geq 2$ . Then  $a^n = a^{p_1} \sqcup\sqcup a^{p_2} \sqcup\sqcup \dots \sqcup\sqcup a^{p_k}$  for every  $p_1 + p_2 + \dots + p_k = n$ , showing that the solution is not unique at all. An extreme case is when  $a^n$  is replaced by  $n$  copies of  $a$ . It appears that this is the only case of ambiguity. More precisely, we show that,

\* Corresponding author.

E-mail address: luc.boasson@liafa.jussieu.fr (L. Boasson).

when powers of a single letter are decomposed into sums of letters (a process we call normalization) then the solution is unique. We call *shuffle root* of  $S$  any normalized solution.

Shuffle factorization has been studied extensively when multiplicities are considered. There are quite precise descriptions of the shuffle algebra  $Q\langle A \rangle$  of polynomials with rational coefficients. A theorem of Radford [5] states that the shuffle algebra  $Q\langle A \rangle$  is freely generated by the Lyndon words (see, e.g., [6, Theorem 6.1]). This means in particular that every finite set  $S$ , viewed as a polynomial with coefficients 0 and 1 is a shuffle of a finite polynomial over Lyndon words, but perhaps with negative coefficients. For instance, since  $a \sqcup b = ab + ba$  and  $ab$  is a Lyndon word (assuming  $a < b$ ), one has  $ba = ab - a \sqcup b$ .

Our approach is different. We consider only ordinary sets of words that are obtained as shuffles of (multi)sets of words, without considering multiplicities in the result, and without considering “negative” multiplicities in the words to be shuffled. The main result is the following.

**Theorem 1.1.** *A finite shuffle set has a unique shuffle root.*

Obviously, this root is effectively computable. In other words, given a finite set  $S$  of nonempty words, there exists at most one effectively computable normalized multiset  $\mathcal{P}$  such that  $S = \sqcup \mathcal{P}$ .

Our method is constructive. Given the set  $S$ , we try to compute some polynomial  $\mathcal{P}$  that is a candidate for the equality  $S = \sqcup \mathcal{P}$ . We show that if the computation fails, the set  $S$  is not a shuffle set. If the computation succeeds, it yields a unique polynomial  $\mathcal{P}$ . If  $S \neq \sqcup \mathcal{P}$ , then  $S$  is not a shuffle set. Otherwise,  $\mathcal{P}$  is the unique polynomial such that  $S = \sqcup \mathcal{P}$ . Thus, our proof gives at most one “candidate” polynomial. Even if this candidate exists, it must be checked that it indeed produces the set  $S$ .

The construction of the candidate polynomial is in two parts. The first is a reduction step that we call separation. The second part is a divide-and-conquer method that strips off parts of the words, and parts in the words to get simpler and simpler sets. A full example of the separation process is given at the end of Section 3. A large example of the second step appears in Section 6.

The separation reduces the computation to the case where all words in  $S$  start with a fixed letter (and all end with this letter, or all end with other letters). The candidate polynomials for the separated sets are computed in the second part. Then they are composed to yield, if the composition succeeds, a candidate that is checked for being the shuffle root.

The initial letter, in the separation process, is called later the pivot. It is used as a marker for decomposing words into blocks. Each block starts with pivot letters, and contains no other pivot letters. It appears that shuffling words can be restricted, with some care, to shuffling sequences of blocks: the blocks play the role of some metaalphabet.

The second part uses two kinds of reductions, called shadow and trim reductions. Both reduce a given separated set as produced in the first step to smaller sets of the same kind. The problem is solved for the two subcases, and their solutions are combined to a solution of the initial set, provided a solution exists. Both reductions operate on blocks. A shadow reduction consists in selecting minimal initial block prefixes, and a trim reduction cuts initial blocks. Considering both reductions is necessary to guarantee that the composed solution is uniquely determined. A commented example is given after Lemma 5.2.

As it will appear, the proof of the theorem also provides an algorithm to test whether a finite set is a shuffle set. This algorithm seems to be much faster than the naive algorithm, however we have no proof for this. Algorithms exist for related problems. Thus, Ref. [7] computes the shuffle of two words without repetition. An algorithm for the computation of the shuffle of a finite set of words, that is efficient in time and in space, has been developed recently by Allauzen [1]. It is based on a clever use of suffix trees. He shows that the shuffle of  $n$  words  $x_1, \dots, x_n$  can be computed in time and space the multinomial

$$O\left(\binom{|x_1| + \dots + |x_n|}{|x_1|, \dots, |x_n|}\right),$$

and that the shuffle set can be stored in linear space, i.e., in space  $O(\text{Card } x_1 \sqcup \dots \sqcup x_n)$ . Ref. [4] gives a parallel algorithm for testing whether a word is in the shuffle of two words. See also [3]. General properties of the shuffle operation are given in [2, Chapter 6].

## 2. Notation

We use standard notation. The empty word is denoted by  $\varepsilon$ . We recall the definition of radix order. Let  $A$  be an ordered alphabet. The *radix order* over  $A^*$  is a total order defined by  $u < v$  if  $|u| < |v|$  or  $|u| = |v|$  and  $u < v$  in the lexicographic order on words (of same length). This order is clearly prefix preserving, i.e.,  $x < xy$  for every nonempty word  $y$ .

*Shuffle*: Let  $A$  be an alphabet. The *shuffle* of two words over  $A$  is defined for words  $x, y$  and letters  $a, b$  by

$$\begin{aligned} x \sqcup \varepsilon &= \varepsilon \sqcup x = x, \\ ax \sqcup by &= a(x \sqcup by) \cup b(ax \sqcup y). \end{aligned}$$

The shuffle is a binary commutative and associative operation. We will be concerned with the shuffle of nonempty words. The shuffle operation is extended to sets as usual by

$$X \sqcup Y = \bigcup_{x \in X, y \in Y} x \sqcup y.$$

We may also consider shuffling all words in a given set. Thus, for a set  $X = \{x_1, x_2, \dots, x_n\}$ , we can define  $\sqcup(X) = x_1 \sqcup x_2 \sqcup \dots \sqcup x_n$ . This *unary* shuffle operation is in fact defined for multisets. Indeed consider for instance

$$a \sqcup ab \sqcup ab = \{ababa, aba^2b, a^2bab, a^2b^2a, a^3b^2\}.$$

This is quite different from  $a \sqcup ab = \{aba, a^2b\}$ . We call *shuffle set* any finite set  $S$  of nonempty words that is the shuffle of a multiset of words.

*Multisets:* A finite *multiset* is a finite set  $\mathcal{P}$  with a multiplicity for each of its elements. It can also be viewed as a *polynomial* with nonnegative coefficients. We denote by  $(\mathcal{P}, w)$  the multiplicity of  $w$  as an element of  $\mathcal{P}$ . It is convenient to write  $\mathcal{P}$  as a polynomial, each word carrying its multiplicity. Thus,  $\mathcal{P} = a + 2ab$  means that  $ab$  appears two times in  $\mathcal{P}$  and  $a$  once.

The *underlying set* or *support* of  $\mathcal{P}$  is the set of words with positive multiplicity. As an example, the support of the multiset  $\mathcal{P} = a + 2ab$  is  $\{a, ab\}$ . We will say that a word  $w$  is an element of the multiset, or *is in* the multiset  $\mathcal{P}$ , if it is an element of the support of  $\mathcal{P}$ .

Observe that, in our problem, the result of shuffling a multiset is viewed as an ordinary set without multiplicities. One could also consider the multiset obtained by shuffling words in  $\mathcal{P}$  by counting the number of times each word is obtained. In the above example, this would give  $2ababa + 4aba^2b + 8a^2bab + 4a^2b^2a + 12a^3b^2$ . However, in this paper, we only consider the ordinary sets produced by shuffling multisets. That is, we forget the coefficients in the result.

The difference  $\mathcal{P} - \mathcal{P}'$  of two multisets  $\mathcal{P}$  and  $\mathcal{P}'$  is defined if  $(\mathcal{P}, w) \geq (\mathcal{P}', w)$  for all words  $w$ . It is then given by  $(\mathcal{P} - \mathcal{P}', w) = (\mathcal{P}, w) - (\mathcal{P}', w)$ .

*Normalization:* A multiset  $\mathcal{P}$  is *proper* if it does not contain the empty word, that is if  $(\mathcal{P}, \varepsilon) = 0$ . A proper multiset  $\mathcal{P}$  is *normal for the letter  $a$*  if the only power of  $a$  that is an element of  $\mathcal{P}$  is  $a$  itself. Thus,  $\mathcal{P} = 37a + 2ab$  is normal for  $a$ , but  $\mathcal{P}' = \varepsilon + a^2 + 7a^5 + 2ab$  is not. A *normal* multiset is a multiset that is proper and normal for all letters.

The *normalization* of a multiset  $\mathcal{P}$  consists in two steps. First, the empty word is removed. Second, the powers of a letter  $a$  that are in  $\mathcal{P}$  are “summed up” to a single monomial  $H * a$ , with  $H = \sum_{n>0} (\mathcal{P}, a^n) * n$ . Thus, for  $\mathcal{P}' = \varepsilon + a^2 + 7a^5 + 2ab$  and for the letter  $a$ , one obtains  $H = 2 + 7 * 5 = 37$ , thus the normalization of  $\mathcal{P}'$  yields  $\mathcal{P} = 37a + 2ab$ . The normalized multiset of  $\mathcal{Q}$  is denoted by  $v(\mathcal{Q})$ .

The reason for doing normalization is the following. The powers of a letter  $a$  contribute a single word in the shuffle obtained by the formula given for  $H$ . Thus,  $\sqcup(a^2 + 7a^5) = a^{37}$ , and  $\sqcup(a^2 + 7a^5 + \mathcal{Q}) = \sqcup(37a + \mathcal{Q})$  for every polynomial  $\mathcal{Q}$ . Therefore, replacing powers of a letter by an appropriate multiplicity of the letter will produce the same shuffle set.

A normalized multiset  $\mathcal{P}$  such that  $S = \sqcup \mathcal{P}$ , is called a *shuffle root* of  $S$ .

### 3. Separation

In this section, we introduce a reduction called separation. This operation reduces a shuffle set and its shuffle root to a simpler pair, with the property that every word in the reduced shuffle root starts with a given letter (and ends with this letter or never ends with this letter).

Two kinds of separations will be defined. The first is applied to sets that are supposed to be shuffle sets, the second one to multisets that are supposed to be their shuffle root. We prove in Lemma 3.3 that the relation between a shuffle set and its shuffle root is preserved through separation. Moreover, Proposition 3.4 shows that the shuffle root of the initial set, if it exists, is determined by the shuffle roots of the separated sets.

A set of words  $S$  is called *a-simple* if  $S \subset aA^* \setminus A^*a$ , that is if all its words start with the letter  $a$ , and none of its words end with  $a$ . A set  $S$  is *a-symmetric* if  $S \subset aA^* \cap A^*a$ , that is if all its words start and end with the letter  $a$ . A set is simple (symmetric) if it is *a-simple* (*a-symmetric*) for some  $a$ . A multiset is simple (symmetric) if its support is simple (symmetric).

For instance, the set  $S = \{a^2bc, acb\}$  is *a-simple*. Clearly, every set  $S \subset A^+$  is a disjoint union of simple sets and symmetric sets.

Let  $B$  be a set of letters. For each word  $w \in A^*$ , we denote by  $\langle B \rangle w$  the longest suffix of  $w$  starting with a letter in  $B$ , if such a suffix exists, the empty set otherwise. Equivalently,  $\langle B \rangle w$  is the word that remains when the longest prefix in  $(A \setminus B)^*$  is removed from  $w$  (and is the empty set when this longest prefix is  $w$  itself). For an intuitive explanation, see also the comment following Eq. (2).

For instance, if  $B = \{b, c\}$ , one has  $\langle B \rangle a^2dba^2c = ba^2c$ . This operation is extended to multisets by linearity. Symmetrically, one defines  $w \langle B \rangle$ .

**Example.** Consider  $\mathcal{P} = a + 2ab$ . Then  $\langle a \rangle \mathcal{P} = \mathcal{P}$ ,  $\langle b \rangle \mathcal{P} = 2b$ ,  $\mathcal{P} \langle a \rangle = 3a$ ,  $\mathcal{P} \langle b \rangle = 2ab$ ,  $\langle \langle a \rangle \mathcal{P} \rangle \langle b \rangle = \langle a \rangle (\mathcal{P} \langle b \rangle) = 2ab$ .

The following lemma will be used in the sequel for two special partitions. In the first partition, each  $A_i$  will consist of a single letter. The second partition will have two sets, one composed of a single letter, the other being the rest of the alphabet.

**Lemma 3.1.** *Let  $A = A_1 \cup A_2 \cup \dots \cup A_n$  be a partition of the alphabet  $A$ . Given normalized multisets  $\mathcal{P}_1, \mathcal{P}_2, \dots, \mathcal{P}_n$ , there exists at most one, effectively computable normalized multiset  $\mathcal{P}$  such that  $\langle A_i \rangle \mathcal{P} = \mathcal{P}_i$  for  $i = 1, \dots, n$ .*

**Proof.** Clearly, if  $\langle A_i \rangle \mathcal{P} = \mathcal{P}_i$  then every word in  $\mathcal{P}_i$  starts with a letter in  $A_i$ . Let  $\ell_i$  be the maximal length of the words in  $\mathcal{P}_i$ , for  $i = 1, \dots, n$ . If  $\ell_i = 0$ , then  $\mathcal{P}_i = 0$ . If  $\ell_i = 1$ , then  $\mathcal{P}_i$  is just a linear combination of letters:  $\mathcal{P}_i = \sum_{a \in A_i} (\mathcal{P}_i, a)a$ . Thus, if  $\ell_i \leq 1$  for  $i = 1, \dots, n$ , then the multiset  $\mathcal{P} = \mathcal{P}_1 + \dots + \mathcal{P}_n$  is a solution, and it is the only solution. Arguing by induction on  $\ell = \max(\ell_1, \dots, \ell_n)$ , consider the polynomial  $\mathcal{Q} = \sum_{i=1}^n \sum_{|x|=\ell} (\mathcal{P}_i, x)x$ . Each word  $x$  in  $\mathcal{Q}$  is also in  $\mathcal{P}$  and conversely, each word

of length  $\ell$  in  $\mathcal{P}$  is in  $\mathcal{Q}$ , because the support of  $\mathcal{P}$  is contained in the union of the supports of the  $P_i$ . Thus  $\mathcal{Q} = \sum_{|x|=\ell} (\mathcal{P}, x)x$ . Consider  $\mathcal{P}' = \mathcal{P} - \mathcal{Q}$  and  $\mathcal{P}'_i = \mathcal{P}_i - \langle A_i \rangle \mathcal{Q}$ . Then  $\langle A_i \rangle \mathcal{P} = \mathcal{P}_i$  if and only if  $\langle A_i \rangle \mathcal{P}' = \mathcal{P}'_i$ , and the lemma is proved by induction.  $\square$

**Example.** Consider the alphabet  $A = \{a, b\}$ . We look for a multiset  $\mathcal{P}$  such that  $\langle a \rangle \mathcal{P} = 2a + ab$  and  $\langle b \rangle \mathcal{P} = b + ba$ . Words of length 2 in  $\mathcal{P}$  are  $ab$  and  $ba$ , so  $\mathcal{P} = ab + ba + \mathcal{P}'$  for some  $\mathcal{P}'$  composed of letters. Next, we compute  $\langle a \rangle \mathcal{P} = ab + a + \langle a \rangle \mathcal{P}'$ , and since  $\langle a \rangle \mathcal{P} = 2a + ab$ , it follows that  $\langle a \rangle \mathcal{P}' = a$ . Similarly,  $\langle b \rangle \mathcal{P} = b + ba + \langle b \rangle \mathcal{P}'$ , whence  $\langle b \rangle \mathcal{P}' = 0$ . Thus,  $\mathcal{P}' = a$  and  $\mathcal{P} = a + ab + ba$ .

As a second example, we look for a multiset  $\mathcal{P}$  such that  $\langle a \rangle \mathcal{P} = a$  and  $\langle b \rangle \mathcal{P} = 2ba$ . The word  $ba$  is the only word of length 2 in  $\mathcal{P}$ . Thus,  $\mathcal{P} = 2ba + \mathcal{P}'$ , and  $a = \langle a \rangle \mathcal{P} = 2a + \langle a \rangle \mathcal{P}'$ . Thus, the multiset we look for does not exist.

Let  $A$  be an ordered alphabet, and let  $B$  be a subset of  $A$ . Given a finite set  $S$ , we consider the set of words over the alphabet  $A \setminus B$  that are prefixes of words in  $S$ . Among these, we choose the word  $p$  that is maximal for the radix order. We denote by  $[B]S$  the set

$$[B]S = p^{-1}S \cap BA^*.$$

Symmetrically, one defines  $S[B]$ . The choice of  $p$  in the definition above is rather arbitrary. It is only necessary, for the propositions below, that  $p$  is a prefix of maximal length.

**Example.** Let  $S = \{ababa, aba^2b, a^2bab, a^2b^2a, a^3b^2\}$ . Then  $[a]S = S$ ,  $[b]S = (a^3)^{-1}S = \{b^2\}$ ,  $S[a] = S(b^2)^{-1} = \{a^3\}$ ,  $S[b] = Sa^{-1} = \{abab, a^2b^2\}$ .

The next proposition shows that, for shuffle sets, this has no consequence.

**Proposition 3.2.** *Let  $S$  be a finite set, let  $B$  be a subset of  $A$ , and let  $p, q$  be words over the alphabet  $A \setminus B$  of maximal length that are prefixes of words in  $S$ . If  $S$  is a shuffle set, then  $p^{-1}S = q^{-1}S$ .*

**Proof.** Assume that  $S = x_1 \sqcup \dots \sqcup x_n$ , and set  $x_i = p_i y_i$ , where  $p_i$  is the longest prefix of  $x_i$  that is in  $(A \setminus B)^*$  and  $y_i \in \varepsilon \cup BA^*$ . Then  $p$  and  $q$  both are in the set  $p_1 \sqcup \dots \sqcup p_n$ , and

$$p^{-1}S = q^{-1}S = y_1 \sqcup \dots \sqcup y_n. \quad \square$$

The lemma does not hold if the set  $S$  is not a shuffle set. Indeed, the set  $S = \{badc, abcd\}$  is not a shuffle set. For  $B = \{c, d\}$ , the prefixes  $ba$  and  $ab$  have maximal length. However,  $(ba)^{-1}S = \{dc\} \neq (ab)^{-1}S = \{cd\}$ .

The next lemma shows how a pair (shuffle set, shuffle root) goes through a separation.

**Lemma 3.3.** *Let  $A = A_1 \cup A_2 \cup \dots \cup A_n$  be a partition of the alphabet  $A$ , and let  $S$  be a finite set over  $A$ . If  $S = \sqcup \mathcal{P}$  for some multiset  $\mathcal{P}$ , then  $[A_i]S = \sqcup \langle A_i \rangle \mathcal{P}$  for each  $i = 1, \dots, n$ .*

**Proof.** As in the previous lemma, we factorize any word  $x$  in  $\mathcal{P}$  into  $x = p_i(x)y_i(x)$ , where  $p_i(x)$  is the longest prefix of  $x$  that is in  $(A \setminus A_i)^*$  and  $y_i(x) \in \varepsilon \cup A_i A^*$ . Then  $[A_i]S = \sqcup \sum (\mathcal{P}, x) y_i(x)$ . Since  $\sum (\mathcal{P}, x) y_i(x) = \langle A_i \rangle \mathcal{P}$ , the lemma is proved.  $\square$

The lemma can also be stated in a more compact form as follows. If  $\mathcal{P}$  is a polynomial and  $B$  is a subset of the alphabet  $A$ , then

$$\sqcup \langle B \rangle \mathcal{P} = [B] \sqcup \mathcal{P} \tag{2}$$

showing how the shuffle operation passes through the  $[ ]$ -operator.

As an explanation of (2), consider three words  $u = pax$ ,  $v = qay$ ,  $w = raz$ , where  $p, q, r$  do not contain the letter  $a$ . Any word in the shuffle of these words has the form  $tas$  where  $t$  does not contain the letter  $a$ . The prefix  $t$  is in the shuffle of some prefixes of  $u, v, w$ , and at least one of the words  $p, q, r$  has been “used” entirely for the construction of  $t$ . If we look for words in the shuffle with the longest possible prefix  $t$ , then  $t$  is in the shuffle  $p \sqcup q \sqcup r$ , and  $as$  is in the shuffle  $ax \sqcup ay \sqcup az$ .

The next proposition is the first special case of partition, concerning a partition into letters.

**Proposition 3.4.** *Let  $S$  be a finite set.*

- (i) *If  $S$  is a shuffle set, then  $[a]S$  is a shuffle set for each letter  $a \in A$ .*
- (ii) *Conversely, if  $[a]S$  is a shuffle set with shuffle root  $\mathcal{P}_a$  for each  $a \in A$ , then there exists at most one, effectively computable multiset  $\mathcal{P}$  such that  $S = \sqcup \mathcal{P}$  and  $\langle a \rangle \mathcal{P} = \mathcal{P}_a$  for  $a \in A$ .*

**Proof.** The first implication is a direct consequence of Lemma 3.3. The converse implication is just Lemma 3.1.  $\square$

It is not difficult to show that, given subsets  $B, C$  of the alphabet, and a set  $S$  of words, the sets  $([B]S)[C]$  and  $[B](S[C])$  are equal. Similarly, for a polynomial  $\mathcal{P}$ , one has  $(\langle B \rangle \mathcal{P}) \langle C \rangle = \langle B \rangle (\mathcal{P} \langle C \rangle)$ . For clarity, we introduce the notation  ${}_B S_C = [B]S[C]$  when  $S$  is a set, and  ${}^B \mathcal{P}^C = \langle B \rangle \mathcal{P} \langle C \rangle$  when  $\mathcal{P}$  is a polynomial.

We will be interested in the sets  ${}_a S_a$  and  ${}_a S_{A \setminus a}$ , and in the polynomials  ${}^a \mathcal{P}^a$  and  ${}^a \mathcal{P}^{A \setminus a}$  for  $a \in A$ .

**Theorem 3.5.** *Let  $S$  be a finite set of words.*

- (i) *If  $S$  is a shuffle set, then the sets  ${}_a S_a$  and  ${}_a S_{A \setminus a}$  ( $a \in A$ ) are shuffle sets.*
- (ii) *Conversely, if the sets  ${}_a S_a$  and  ${}_a S_{A \setminus a}$  ( $a \in A$ ) are shuffle sets, and if polynomials  $\mathcal{Q}_a$  and  $\mathcal{R}_a$  are given such that  ${}_a S_a = \sqcup \mathcal{Q}_a$  and  ${}_a S_{A \setminus a} = \sqcup \mathcal{R}_a$  for  $a \in A$ , then there exists at most one effectively computable polynomial  $\mathcal{P}$  such that  $S = \sqcup \mathcal{P}$  and  ${}^a \mathcal{P}^a = \mathcal{Q}_a$  and  ${}^a \mathcal{P}^{A \setminus a} = \mathcal{R}_a$ .  $\square$*

In view of this result, the computation of a shuffle root of a set can be reduced to the cases we called  $a$ -simple and  $a$ -symmetric. In other words, a given set  $S$  for which a shuffle root is looked for, is first decomposed into simple and symmetric sets  ${}_a S_a$  and  ${}_a S_{A \setminus a}$  ( $a \in A$ ). For each of these sets, a shuffle root is computed, provided it exists. There is at most one polynomial that can be constructed by composing these roots. The initial set  $S$  is a shuffle set only if the composed polynomial exists. Even then, the set  $S$  might not be a shuffle set. For a detailed example, see Section 6. For further reference, we state this fact as a corollary.

**Corollary 3.6.** *Let  $S$  be a finite set of nonempty words. If, for each  $a \in A$ , the  $a$ -simple set  ${}_a S_a$  and  $a$ -symmetric set  ${}_a S_{A \setminus a}$  ( $a \in A$ ) is a shuffle set and has a unique shuffle root, then the set  $S$  has at most one effectively computable shuffle root.  $\square$*

The proof of Theorem 1.1 thus reduces to the case of simple and symmetric sets. We finish this section with an example illustrative of the separation process. See also Section 6.

**Example.** Let  $S$  be the set

$$\begin{array}{l} aabab \ aabba \ ababa \ abaab \ abbaa \\ baaab \ baaba \ babaa. \end{array}$$

Then

$$[a]S = \{aaab, aaba, abaa\}, \quad [b]S = \{bab, bba\}$$

and

$$[a]S[a] = \{a^3\}, \quad [a]S[b] = \{ab\}, \quad [b]S[a] = \{ba\}, \quad [b]S[b] = \{b^2\}.$$

Since these four sets are singletons, they are their own shuffles, and after normalization, we get

$$[a]S[a] = \sqcup\sqcup 3a, \quad [a]S[b] = \sqcup\sqcup ab, \quad [b]S[a] = \sqcup\sqcup ba, \quad [b]S[b] = \sqcup\sqcup 2b.$$

If  $[a]S = \sqcup\sqcup \mathcal{P}$  for a polynomial  $\mathcal{P}$ , then  $\mathcal{P} = \alpha ab + \beta a$  for some nonnegative integers  $\alpha, \beta$ . From  $\mathcal{P}\langle a \rangle = (\alpha + \beta)a = 3a$  and  $\mathcal{P}\langle b \rangle = \alpha ab = ab$  Content-Length: 22322.

From  $\mathcal{P}\langle a \rangle = (\alpha + \beta)a = 3a$  and  $\mathcal{P}\langle b \rangle = \alpha ab = ab$  we get  $\mathcal{P} = ab + 2a$ .

Thus the only possibility is  $[a]S = \sqcup\sqcup (2a + ab)$ . Similarly, one must have  $[b]S = \sqcup\sqcup (b + ba)$ . Next, if  $S = \sqcup\sqcup \mathcal{Q}$ , then

$$\mathcal{Q} = \alpha a + \beta ab + \gamma b + \delta ba$$

for some nonnegative integers  $\alpha, \beta, \gamma, \delta$ , and since

$$\langle a \rangle \mathcal{Q} = (\alpha + \delta)a + \beta ab = 2a + ab, \quad \langle b \rangle \mathcal{Q} = (\beta + \gamma)b + \delta ba = b + ba$$

we get  $\alpha = \beta = \delta = 1, \gamma = 0$  and  $\mathcal{Q} = a + ab + ba$ . A final check shows that indeed  $S = a \sqcup\sqcup ab \sqcup\sqcup ba$ . Observe that in the separation process, only five of the words in  $S$

were used. This means that if we had omitted any of the unused words *ababa*, *abaab* or *abbaa*, the candidate polynomial  $\mathcal{Q}$  would have been the same. This shows that the final check is mandatory.

#### 4. Blocks

In this section, we consider a canonical decomposition of words induced by a distinguished letter  $a$ , assuming all words start with this letter. An intricate-order relation will be defined that allows to order blocks. The main observation is that minimal words (for the radix order) in a shuffle set can be computed by considering minimal words for the block order in the shuffle root. These properties will be used in the next section.

*Pivot and blocks:* We fix a distinguished letter  $a$  to be the *pivot* and set  $B = A \setminus \{a\}$ . A *block* is any word in  $a^+B^*$ , a *regular block* is a word in  $a^+B^+$ . A block in  $a^+$  is *irregular*. Any word  $w$  in  $aA^*$  has a unique factorization into blocks. Only the last block of this factorization may be irregular. The *initial block* of  $x \in aA^*$  is the first block in its factorization. The *block degree*  $\text{deg}(x)$  of a word  $x \in aA^*$  is the number of blocks in its factorization.

*Order:* We consider any total order relation  $<$  on  $B^*$  that is compatible with the prefix ordering, i.e., such that  $x < xy$  for  $x \in B^*$ ,  $y \in B^+$ . The radix order is an example of such an order. We define an order on blocks called *block order*, and also denoted by  $<$ , by the following requirements. First, irregular blocks are ordered by their exponent:  $a^p < a^n$  iff  $p < n$ . Next, any irregular block is greater than any regular block:  $a^p u < a^n$  for any  $n, p$  and every nonempty word  $u$  over  $B$ . Finally, regular blocks are ordered first by the exponent of the pivot and, if these are equal, by the *opposite* of the order on  $B^*$ : for  $u, v \in B^+$ ,

$$a^p u < a^n v \Leftrightarrow p < n \quad \text{or} \quad (n = p \text{ and } u > v).$$

Observe the inversion of the order on  $u$  and  $v$ .

The block order is extended on  $aA^*$  by the lexicographic extension of the order on blocks. For this, consider

$$x = \beta_1 \beta_2 \cdots \beta_n, \quad x' = \beta'_1 \beta'_2 \cdots \beta'_{n'}.$$

We say that  $x$  is a proper *block prefix* of  $x'$  if  $n < n'$  and  $\beta_i = \beta'_i$  for  $1 \leq i \leq n$ . We next define  $x \sqsubset x'$  if and only if there is an integer  $k \leq n, n'$  such that  $\beta_i = \beta'_i$  for  $1 \leq i < k$  and  $\beta_k < \beta'_k$ , and we set  $x < x'$  if and only if  $x \sqsubset x'$  or  $x$  is a proper block prefix of  $x'$ .

For blocks  $\beta$  and  $\beta'$ , one has  $\beta < \beta'$  iff  $\beta \sqsubset \beta'$ .

**Example.** Consider the block order induced by the radix order on  $B = \{b, c\}$  with  $b < c$ , and consider words  $x = a^2 b^5 c^5 b a^2 b^3 c^5 b$ ,  $x' = a^2 b^3 c^3 b$ ,  $x'' = a^2 b^5 c^5 b a$ . We get block factorizations  $x = \beta_1 \beta_2$ ,  $x' = \beta'_1$ ,  $x'' = \beta_1 \gamma$  with  $\beta_1 = a^2 b^5 c^5 b$ ,  $\beta_2 = a^2 b^3 c^5 b$ ,  $\beta'_1 = a^2 b^3 c^3 b$ ,  $\gamma = a$ . To show that  $x < x'$ , we compare  $\beta_1$  and  $\beta'_1$ . They have the same initial power of

the letter  $a$ , namely  $a^2$ . Since  $|\beta_1| > |\beta'_1|$ , we get  $\beta_1 < \beta'_1$  and consequently  $x < x'$ . Next consider  $x$  and  $x''$ . The word  $x''$  is a prefix of  $x$  as a word of  $A^*$ . However,  $\beta_2 < \gamma$  and therefore  $x < x''$ .

*Properties of the order:* The order relation over blocks helps in reducing (frequently) shuffle of words to rearrangements of blocks by inspection of special representatives. Given a set  $X \subset aA^*$ , we denote by  $\min(X)$  the smallest element in  $X$  for the block order. We shall see that if  $x = \beta_1 \cdots \beta_n$  and  $y = \gamma_1 \cdots \gamma_m$ , then  $\min(x \sqcup y)$  is the product of  $\beta_1, \dots, \beta_n, \gamma_1, \dots, \gamma_m$  in some order.

First, given blocks  $\beta, \gamma$ , one has

$$\beta \leq \gamma \Rightarrow \min(\beta \sqcup \gamma) = \beta\gamma.$$

Observe that  $\beta\gamma$  might be a single block if  $\beta, \gamma$  are both irregular. Next, if  $x$  has the block decomposition  $x = \beta_1\beta_2 \cdots \beta_n$  for some  $n \geq 2$ , then

$$\min(\beta \sqcup x) = \beta_1\beta_2 \cdots \beta_i\beta\beta_{i+1} \cdots \beta_n,$$

where  $i$  is the greatest index such that  $\beta_j < \beta$  for  $j = 1, \dots, i$ . For instance, if  $\beta = a^2b^2$  and  $x = ab^3 \cdot a^2b^3 \cdot a^2b$ , then  $\min(\beta \sqcup x) = ab^3 \cdot a^2b^3 \cdot a^2b^2 \cdot a^2b$ . The relation extends to words in  $aA^*$  as follows. Let  $x = \beta x'$ ,  $y = \gamma y'$  for (regular) blocks  $\beta, \gamma$  and words  $x', y' \in aA^*$ . Then

$$\min(x \sqcup y) = \begin{cases} \beta \min(x' \sqcup y) & \text{if } \beta < \gamma, \\ \beta \min(\min(x' \sqcup y), \min(x \sqcup y')) & \text{if } \beta = \gamma. \end{cases}$$

Indeed, assume  $\beta \leq \gamma$  and set  $\beta = a^n u$  and  $\gamma = a^p v$  ( $u, v \neq \varepsilon$ ). The word  $\min(x \sqcup y)$  starts with the smallest possible numbers of  $a$ 's, that is by  $a^n$ . This group of  $a$ 's is followed by the longest possible word containing no  $a$ . If  $n < p$ , this must be  $u$ . If  $n = p$ , it is also  $u$  if  $\beta < \gamma$ , and if  $\beta = \gamma$  then it is  $u = v$ . Thus, in all cases,  $\min(x \sqcup y)$  starts with  $\beta$ . It is followed by the smallest possible word. This is  $\min(x' \sqcup y)$  if  $\beta < \gamma$ . If  $\beta = \gamma$ , the smallest possible word is either  $\min(x' \sqcup y)$  if the initial block  $\beta$  of  $\min(x \sqcup y)$  is the initial block of  $x$ , or  $\min(x \sqcup y')$  if the block  $\beta$  is the initial block of  $y$ . This proves the formula. The three formulas allow to compute the minimal word in the shuffle of two words in  $aA^*$ .

Next, for blocks  $\beta, \gamma$  and the word  $x \in aA^*$ ,

$$\beta < \gamma \Rightarrow x\beta \sqsubset x\gamma$$

and for blocks  $\beta, \gamma$  and a words  $x, y \in aA^*$ ,

$$\beta < \gamma \quad \text{and} \quad \gamma \text{ regular} \Rightarrow \beta x \sqsubset \gamma y.$$

The last relation does not hold if  $\gamma$  is irregular. Consider for instance  $\beta = a$ ,  $\gamma = a^2$  and take  $x = a^3$  and  $y = a$ . Then  $\beta < \gamma$  and  $a^4 = \beta x \sqsubset a^3$ .

Next, we show that for  $x, y, y' \in aA^*$ ,

$$y \sqsubset y' \Rightarrow \min(x \sqcup y) \sqsubset \min(x \sqcup y').$$

Indeed, set  $y = p\beta s$ ,  $y' = p\gamma s'$  for blocks  $\beta < \gamma$  and words  $p, s, s' \in aA^* \cup \{\varepsilon\}$ . Consider  $w' = \min(x \sqcup y')$ . There is a factorization  $x = x_1 x_2$  of  $x$  into products of blocks such that  $w' = u\gamma v'$  for words  $u \in x_1 \sqcup p$ ,  $v' \in x_2 \sqcup s'$ . Consider  $w = u\beta v$ , for  $v \in x_2 \sqcup s$ . If  $\gamma$  is regular, then  $w \sqsubset w'$ . If  $\gamma$  is irregular, then  $s'$  is the empty word and  $v' = x_2$ . Thus,  $x_2$  is an irregular block, because  $\gamma x_2 = \gamma v'$  is the last block of  $w'$ . Since  $x_2 \gamma = \gamma x_2$ , we may assume  $x = x_1$ ,  $x_2 = \varepsilon$  in the factorization of  $x$ . Thus,  $w = u\beta$  and again  $w \sqsubset w'$ . Thus,  $\min(x \sqcup y) \leq w \sqsubset w' = \min(x \sqcup y')$ .

**Lemma 4.1.** *Let  $Y$  be a subset of  $aA^*$ . If no word in  $Y$  is a proper block prefix of another word in  $Y$ , then for any word  $x$  in  $aA^*$ ,*

$$\min(x \sqcup Y) = \min(x \sqcup \min(Y)).$$

**Proof.** By definition,  $\min(x \sqcup Y) = \min_{y \in Y} (x \sqcup y)$ . Since  $y \sqsubset y'$  implies that  $\min(x \sqcup y) \sqsubset \min(x \sqcup y')$ , the result follows by induction on the size of  $Y$ .  $\square$

**Example.** Consider  $x = a^2 bab$  and  $Y = \{a^2 bab, a^2 ba^2 b, a^2 ba^3 b, a^4 b\}$ . Then  $\min(x \sqcup a^2 bab) = a^2 bab a^2 bab$ ,  $\min(x \sqcup a^2 ba^2 b) = a^2 baba^2 ba^2 b$ ,  $\min(x \sqcup a^2 ba^3 b) = a^2 baba^2 ba^3 b$ ,  $\min(x \sqcup a^4 b) = a^2 baba^4 b$ , and the smallest of these words is the first one, which is  $\min(x \sqcup \min(Y))$ .

The lemma implies that

$$\min(x \sqcup y \sqcup z) = \min(x \sqcup \min(y \sqcup z))$$

since  $x \sqcup y \sqcup z = x \sqcup Y$  with  $Y = y \sqcup z$ . This formula, together with the formulas given earlier, allow to compute the minimal word in the shuffle of any polynomial. Observe that, up to irregular blocks, this minimal word is obtained as some permutation of the blocks composing the words in the support, each block appearing with the appropriate multiplicity.

*Parikh vector:* The *Parikh vector*  $\text{par}(x)$  of a word  $x$  is defined as usual. Let  $k$  be the size of the alphabet  $A$ , and for each letter  $c \in A$  let  $|x|_c$  be the number of  $c$ 's in  $x$ . Then  $\text{par}(x) = (|x|_c)_{c \in A}$ . As an example,  $\text{par}(a^2 bca^3 ba^2) = (7, 2, 1)$ . The Parikh vector  $\text{par}(\mathcal{P})$  of a polynomial  $\mathcal{P}$  is defined by linearity:

$$\text{par}(\mathcal{P}) = \sum_{x \in P} (\mathcal{P}, x) \text{par}(x)$$

For instance,  $\text{par}(2a^2 bab + a^2 ba^2 b + a^2 ba^3 b + 2a^4 b) = (23, 10)$ .

If  $S = \sqcup (\mathcal{P})$ , then all words in  $S$  have the same Parikh vector, and this Parikh vector is  $\text{par}(\mathcal{P})$ .

Vectors as ordered as usual. For two vectors  $p, p'$ , we set  $p \leq p'$  iff  $p_c \leq p'_c$  for all letters  $c$ . We will use this order for proof by induction.

## 5. Two reductions

We assume that the pivot is the letter  $a$ . Given a set of words  $S$  in  $aA^*$ , the *minimal initial block* of  $S$  is the smallest of the initial blocks of the words in  $S$ . It is denoted  $\mu_S$ . The minimal initial block of a polynomial  $\mathcal{P}$  is the minimal initial block of its support. It is denoted  $\mu_{\mathcal{P}}$ .

### 5.1. Shadow reduction

Let  $\beta$  be a block. Let  $x$  be a word in  $aA^*$  and let  $x = \beta_1\beta_2 \cdots \beta_n$  be its block decomposition. We define a word  $\pi_{\beta}(x)$  by

$$\pi_{\beta}(x) = \beta_1\beta_2 \cdots \beta_i,$$

where  $i$  is the smallest index such that  $\beta_{i+1} > \beta$ . Of course,  $\pi_{\beta}(x) = \varepsilon$  if  $\beta_1 > \beta$ , and  $\pi_{\beta}(x) = x$  if  $\beta_1, \beta_2, \dots, \beta_n \leq \beta$ . This construction will be used in two slightly different ways.

Let  $y = \gamma_1 \cdots \gamma_m$  be a second word in  $aA^*$ . Set  $x = \pi_{\beta}(x)x'$  and  $y = \pi_{\beta}(y)y'$ . Then  $\min(x \sqcup \sqcup y)$  starts with  $\min(\pi_{\beta}(x) \sqcup \sqcup \pi_{\beta}(y))$  and more precisely

$$\min(x \sqcup \sqcup y) = \min(\pi_{\beta}(x) \sqcup \sqcup \pi_{\beta}(y)) \min(x' \sqcup \sqcup y').$$

**Example.** Let  $x = a^2b^3 \cdot a^2b^2 \cdot a^3b$  and  $y = a^2b^2 \cdot ab \cdot a^2b$ . For  $\beta = a^2b^2$ , one gets  $\pi_{\beta}(x) = a^2b^3 \cdot a^2b^2$  and  $\pi_{\beta}(y) = a^2b^2 \cdot ab$ . Next  $\min(x \sqcup \sqcup y) = a^2b^3 \cdot a^2b^2 \cdot ab \cdot a^2b^2 \cdot a^2b \cdot a^3b$ . Also,  $\min(a^2b^3 \cdot a^2b^2 \sqcup \sqcup a^2b^2 \cdot ab) = a^2b^3 \cdot a^2b^2 \cdot ab \cdot a^2b^2$ .

Let  $S$  be a finite subset of  $aA^*$ , and let  $\mu = \mu_S$  be the minimal initial block of  $S$ . Let  $s = \beta_1\beta_2 \cdots \beta_n$  be the smallest word in  $S$  for the block order. Clearly  $\beta_1 = \mu$ . We define

$$\hat{S} = (\pi_{\mu}(s))^{-1}S \cap aA^*.$$

This set of words is computed in two steps. First, only words starting with  $\pi_{\mu}(s)$  are kept. Next, this prefix is removed from each of the remaining words, and the empty word, if any, is discarded. As a first example, consider the set

$$S = \{ab \cdot ab \cdot a, ab \cdot aab, aab \cdot ab, aabb \cdot a, aaabb\}.$$

Here  $s = ababa$  and  $\mu = \mu_S = ab$ ,  $\pi_{\mu}(s) = abab$  and consequently  $\hat{S} = \{a\}$ . Other examples will be given below.

Consider next a polynomial  $\mathcal{P}$ . Recall that  $\mu = \mu_{\mathcal{P}}$  is the smallest of the initial blocks of the words in the support of  $\mathcal{P}$ . Define

$$\mathcal{P}' = \sum (\mathcal{P}, x)(\pi_{\mu}(x)^{-1}x) \quad \text{and} \quad \hat{\mathcal{P}} = v(\mathcal{P}')$$

Here, for each word in the support of  $x$ , the prefix  $\pi_{\mu}(x)$  is computed and removed. The resulting polynomial is then normalized. As a first example, consider  $\mathcal{P} = a + 2ab$ . Here  $\mu_{\mathcal{P}} = ab$ , and  $\mathcal{P}' = 2\varepsilon + a$ , whence  $\hat{\mathcal{P}} = a$ .

The operations  $S \mapsto \hat{S}$  and  $\mathcal{P} \mapsto \hat{\mathcal{P}}$  are called *shadow reductions*. The first is the shadow reduction of a set, the second the shadow reduction of a polynomial.

**Example 1.** Let

$$\mathcal{P} = 2a^2bab + a^2ba^2b + a^2ba^3b + 2a^4b.$$

The initial blocks are  $a^2b$  and  $a^4b$ , and the first one is smaller, thus  $\mu = \mu_{\mathcal{P}} = a^2b$ . We get  $\pi_{\mu}(a^2bab) = a^2bab$ ,  $\pi_{\mu}(a^2ba^2b) = a^2ba^2b$ ,  $\pi_{\mu}(a^2ba^3b) = a^2b$ ,  $\pi_{\mu}(a^4b) = \varepsilon$ , whence  $\mathcal{P}' = 2\varepsilon + \varepsilon + a^3b + 2a^4b$  and  $\hat{\mathcal{P}} = a^3b + 2a^4b$ . Consider the set

$$S = \sqcup\sqcup(\mathcal{P})$$

The smallest word  $s$  in  $S$  for the block order is  $s = a^2baba^2bab(a^2b)^3a^3b(a^4b)^2$ . Here  $\pi_{\mu_S}(s) = a^2baba^2bab(a^2b)^3$ . It is not difficult to check that  $\mu_{\mathcal{P}} = \mu_S$ . It appears moreover that  $\hat{S} = \sqcup\sqcup(\hat{\mathcal{P}})$ . This is a general fact, stated in Lemma 5.1 below.

**Example 2.** This is a variation of the previous example. Let

$$\mathcal{Q} = 2a + 2a^2baba + a^2ba^2ba^2 + 2a^4ba.$$

The initial blocks are  $a$ ,  $a^2b$  and  $a^4b$ , and the second one is the smallest. Thus  $\mu_{\mathcal{Q}} = a^2b$ . We get  $\mathcal{Q}' = 2a + 2a + a^2 + 2a^4ba$ . Normalization groups  $2a + 2a + a^2$  into  $6a$ , whence  $\hat{\mathcal{Q}} = 6a + 2a^4ba$ . Consider  $T = \sqcup\sqcup(\mathcal{Q})$ . The minimal word  $t$  in  $T$  for the block order is  $t = a^2b \cdot ab \cdot a^2b \cdot ab \cdot a^2b \cdot a^4b \cdot a^4b \cdot a^8$ . We get  $\pi_{\mu_T}(t) = a^2baba^2baba^2ba^2b$ , and  $\hat{T} = \pi_{\mu_T}(t)^{-1}T = \sqcup\sqcup(2a + 2a + a^2 + 2a^4b) = \sqcup\sqcup(\hat{\mathcal{Q}})$ .

**Lemma 5.1.** *If  $S = \sqcup\sqcup(\mathcal{P})$ , then  $\hat{S} = \sqcup\sqcup(\hat{\mathcal{P}})$ .*

**Proof.** Let  $s = \mu\beta_2 \cdots \beta_n$  be the smallest word in  $S = \sqcup\sqcup(\mathcal{P})$ . Then  $\mu$  is the smallest of the initial blocks of words in the support  $P$  of  $\mathcal{P}$ . Thus  $\mu = \mu_{\mathcal{P}} = \mu_S$ .

Let  $Q$  be the set of words in  $P$  with initial block  $\mu$ . Set  $\mathcal{Q} = \sum_{q \in Q} (\mathcal{P}, q)q$  and define  $\mathcal{R} = \mathcal{P} - \mathcal{Q}$ . Each word  $q$  in  $Q$  is written as  $q = \pi_{\mu}(q)q'$  with  $\pi_{\mu}(q)$  nonempty. Then  $\hat{\mathcal{P}} = \mathcal{R} + \mathcal{T}$ , with

$$\mathcal{T} = \sum_{q \in Q, q' \neq \varepsilon} (\mathcal{P}, q)q'.$$

Next  $\pi_{\mu}(s) = \min(\sqcup\sqcup(\sum_{q \in Q} (\mathcal{P}, q)\pi_{\mu}(q)))$ . By the discussions above, this word is a product of blocks. It follows that  $\hat{S} = \sqcup\sqcup(\mathcal{R} + \mathcal{T})$ , that is  $\hat{S} = \sqcup\sqcup(\hat{\mathcal{P}})$ .  $\square$

Observe that in Lemma 5.1, the Parikh vector of  $\hat{\mathcal{P}}$  is strictly less than the Parikh vector of  $\mathcal{P}$ .

### 5.2. Trim reduction

For any word  $x = \beta_1\beta_2 \cdots \beta_n$ , we define

$$\bar{x} = \beta_2 \cdots \beta_n.$$

Given a polynomial  $\mathcal{P}$ , we define  $\bar{\mathcal{P}}$  by linearity and normalization as

$$\bar{\mathcal{P}} = v(\sum (\mathcal{P}, x)\bar{x}).$$

Observe that words composed of a single block disappear. As a first example, consider  $\mathcal{P} = a + 2ab$ . Here  $\bar{\mathcal{P}} = v(\varepsilon + 2\varepsilon) = 0$ .

Given a set  $S \subset aA^*$ , the *leading block*  $\lambda_S$  of  $S$  is one of the initial blocks of words of  $S$ . It is the initial block  $a^n u$  of  $S$  with  $n$  maximal and  $|u|$  maximal among the initial blocks with exponent  $n$ . Observe that this choice is not related to the block order we defined previously. Define

$$\bar{S} = \lambda_S^{-1} S \cap aA^*.$$

The operations  $S \mapsto \bar{S}$  and  $\mathcal{P} \mapsto \bar{\mathcal{P}}$  are called *trim reductions*. Note that these reductions strictly again decrease the corresponding Parikh vectors.

**Example 1 (continued).** Let

$$\mathcal{P} = 2a^2bab + a^2ba^2b + a^2ba^3b + 2a^4b.$$

We get

$$\bar{\mathcal{P}} = 2ab + a^2b + a^3b.$$

Next, we consider  $S = \sqcup\sqcup(\mathcal{P})$  without computing it explicitly. The maximal number of initial  $a$ 's in words in  $S$  is obtained by putting all  $a$ 's of the initial blocks of  $\mathcal{P}$  at the beginning of the words. This gives  $a^{16}$ . Next, we look for a longest word over  $A \setminus a$  that can follow this power of  $a$ . This is obtained here by concatenating all  $b$ 's occurring in all initial blocks, that is  $b^6$ . Thus, the block we look for is  $\lambda_S = a^{16}b^6$ . Clearly,  $\bar{S} = \lambda_S^{-1} S$  is the shuffle of what remains of  $\mathcal{P}$ , that is  $\bar{S} = \sqcup\sqcup(\bar{\mathcal{P}})$ .

**Example 2 (continued).** Let

$$\mathcal{Q} = 2a + 2a^2baba + a^2ba^2ba^2 + 2a^4ba.$$

We get  $\bar{\mathcal{Q}} = 2aba + a^2ba^2 + 2a$ . The first term  $2a$  in  $\bar{\mathcal{Q}}$  disappears, and the last term in  $\bar{\mathcal{Q}}$  comes from  $2a^4ba = 2a$ . Consider  $T = \sqcup\sqcup(\bar{\mathcal{Q}})$ . The maximal number of initial  $a$ 's in words in  $T$  is  $a^{16}$ , and the maximal number of  $b$ 's that can follow the power if  $a$  is  $b^5$ . Thus,  $\lambda_T = a^{16}b^5$  and  $\lambda_T^{-1} T = \sqcup\sqcup(2aba + a^2ba^2 + 2a)$ .

**Lemma 5.2.** *If  $S = \sqcup\sqcup(\mathcal{P})$ , then  $\bar{S} = \sqcup\sqcup(\bar{\mathcal{P}})$ .*

**Proof.** Clear from the previous examples.  $\square$

Observe that in Lemma 5.2, the Parikh vector of  $\bar{\mathcal{P}}$  is strictly less than the Parikh vector of  $\mathcal{P}$ .

**Example.** The aim of this example is to illustrate the two reductions introduced earlier and to give a feeling why these will lead to an effective procedure.

Consider three words  $p = a^2b^2ab^3a^3b$ ,  $q = a^2b^2ab^3a^4b^2$  and  $r = a^3ba^2b^2$  and let  $S = p \sqcup q \sqcup r$  be the shuffle set generated by these words.

The maximum number of letters  $a$  at the beginning of a word in  $S$  is clearly  $2 + 2 + 3 = 7$ . Among the words in  $S$  starting with  $a^7$ , those which have the maximum number of  $b$  following  $a^7$  will clearly begin with  $b^{2+2+1} = b^5$ . In order to get these words in  $S$ , we have to use the three initial blocks of the words  $p, q, r$ . Hence, the set  $(a^7b^5)^{-1}S$  is a simpler shuffle set. More precisely, setting  $p' = ab^3a^3b$ ,  $q' = ab^3a^4b^2$  and  $r' = a^2b^2$ , we have  $(a^7b^5)^{-1}S = p' \sqcup q' \sqcup r'$ . The set  $(a^7b^5)^{-1}S$  is precisely what has been defined as  $\hat{S}$ , and the polynomial  $p' + q' + r'$  is the trimmed reduction of  $p + q + r$ .

Now, we look for words in  $S$  beginning with the minimum number of  $a$ . These are the words beginning with  $a^2$ . Among these words, those which have the maximum number of  $b$  following  $a^2$  will clearly begin with  $a^2b^2$ . These are the words in  $S$  beginning with the smallest initial block. However, the set  $(a^2b^2)^{-1}S$  is not a shuffle set. It is the union of two shuffle sets, obtained by removing the block  $a^2b^2$  either from word  $p$  or from word  $q$ . In other words, one of the shuffle sets is  $p' \sqcup q \sqcup r$  and the other is  $p \sqcup q' \sqcup r$ .

Let us repeat the process of looking for “minimal” blocks of words in  $S$  until a block appears that is greater than  $a^2b^2$ . We have to repeat the process four times, giving rise to the word  $h = a^2b^2ab^3a^2b^2ab^3$  composed of four blocks. Then the quotient  $h^{-1}S$  is a shuffle set: indeed, the word  $h$  has used all blocks less or equal to  $a^2b^2$  which is the minimal initial block in the root of  $S$ ; so the quotient is exactly the set obtained by shuffling the remaining suffixes of the original words, that is  $h^{-1}S = a^3b \sqcup a^4b^2 \sqcup r$ . With the previous notation,  $h = \pi_\mu(s)$ , where  $\mu = a^2b^2$  and  $s$  is the smallest word in  $S$ . The quotient is precisely what has been defined as  $\hat{\hat{S}}$ .

Both of these operations lead to a simpler shuffle set. The rest of the proof follows from the fact that, given the solutions of each subproblem, there is at most one (computable) solution of the original problem.

*Degree:* The *block degree*  $\text{deg}(x)$  of a word  $x \in aA^*$  was defined to be the number of blocks in its decomposition. For a polynomial  $\mathcal{P}$ , the block degree  $\text{deg}(\mathcal{P})$  is the maximum of the degrees of the words in its support. As an example, the polynomial of Example 1 has degree 2, and the polynomial of Example 2 has degree 3. It will be convenient to denote  $\mathcal{P}_k$  the *homogeneous component* of degree  $k$  of  $\mathcal{P}$ , that is  $\mathcal{P}_k = \sum_{\text{deg}(x)=k} (\mathcal{P}, x)x$ . A normalized polynomial  $\mathcal{P}$  of degree  $d$  then writes as  $\mathcal{P} = \mathcal{P}_1 + \mathcal{P}_2 + \dots + \mathcal{P}_d$ .

A polynomial  $\mathcal{P}$  is *simple* if no word in its support contains an irregular block, it is *symmetric* if every word in its support contains (and thus ends with) an irregular block. The polynomial in Example 1 is simple, the polynomial in Example 2 is symmetric, the polynomial  $a + 2ab$  is neither simple nor symmetric. The shadow polynomial and the trim polynomial of a simple (symmetric) polynomial are again simple (symmetric).

For any polynomial  $\mathcal{P}$  of degree  $d \geq 1$ , the polynomial  $\hat{\mathcal{P}}$  has degree  $d$  or less, and the trim polynomial  $\bar{\mathcal{P}}$  has degree  $d - 1$ .

**Proposition 5.3.** *Given two simple normalized polynomials  $\mathcal{H}$  and  $\mathcal{T}$ , a vector  $p$ , and a block  $\mu$ , there is at most one simple normalized effectively computable (simple) polynomial  $\mathcal{P}$  such that  $\hat{\mathcal{P}} = \mathcal{H}$ ,  $\bar{\mathcal{P}} = \mathcal{T}$ ,  $\mu_{\mathcal{P}} = \mu$  and  $\text{par}(\mathcal{P}) = p$ .*

**Proof.** By induction on the degree of  $\mathcal{T}$ . If  $\text{deg}(\mathcal{T}) = 0$ , then  $\mathcal{T} = 0$  because  $\mathcal{T}$  is proper. If  $\mathcal{P}$  exists and  $\mathcal{P} \neq 0$ , then  $\text{deg}(\mathcal{P}) = 1$ . Set  $\mathcal{P} = h\mu + \mathcal{Q}$  with  $h = (\mathcal{P}, \mu)$ . Since  $\hat{\mathcal{P}} = \mathcal{H}$ , the support of  $\mathcal{P}$  does not contain an irregular block. Since  $\mu = \mu_{\mathcal{P}}$ , all blocks in  $\mathcal{Q}$  are greater than  $\mu$ , and consequently  $\hat{\mathcal{P}} = \mathcal{Q}$ , hence  $\mathcal{Q} = \mathcal{H}$ . Thus, if  $\mathcal{P}$  exists, it is of the form

$$\mathcal{P} = h * \mu + \mathcal{H}.$$

Finally, since  $p = h \text{par}(\mu) + \text{par} \mathcal{H}$ , the number  $h$  can be computed. If it is a positive integer, then  $\mathcal{P}$  is uniquely determined, otherwise the polynomial  $\mathcal{P}$  does not exist.

Assume next that  $\text{deg}(\mathcal{T}) = d - 1 \geq 1$ . If a polynomial  $\mathcal{P}$  exists with  $\bar{\mathcal{P}} = \mathcal{T}$  and  $\hat{\mathcal{P}} = \mathcal{H}$ , then  $\text{deg}(\mathcal{P}) = d = 1 + \text{deg}(\mathcal{T})$  and  $\text{deg}(\mathcal{P}) \geq \text{deg}(\mathcal{H})$ . Write  $\mathcal{P} = \mathcal{P}' + \mathcal{P}_d$ , with  $\mathcal{P}' = \mathcal{P}_1 + \dots + \mathcal{P}_{d-1}$ . Set  $\mathcal{P}' = \mu\mathcal{C} + \mathcal{N}$ ,  $\mathcal{P}_d = \mu\mathcal{D} + \mathcal{M}$ , such that the initial blocks of the words of the supports of  $\mathcal{N}$  and  $\mathcal{M}$  are greater than  $\mu$ . Also,  $\mathcal{C}$  has degree at most  $d - 2$ ,  $\mathcal{D}$  has degree of degree  $d - 1$  and  $\mathcal{M}$  has degree  $d$ . From  $\mathcal{P} = \mu\mathcal{C} + \mathcal{N} + \mu\mathcal{D} + \mathcal{M}$ , it follows that  $\mathcal{T} = \bar{\mathcal{P}} = \mathcal{C} + \bar{\mathcal{N}} + \mathcal{D} + \bar{\mathcal{M}}$  and  $\mathcal{H} = \hat{\mathcal{P}} = \hat{\mathcal{C}} + \mathcal{N} + \hat{\mathcal{D}} + \mathcal{M}$ . Thus  $\mathcal{H}_d = \mathcal{M}$ , and  $\mathcal{T}_{d-1} = \mathcal{D} + \bar{\mathcal{M}}$ . This determines  $\mathcal{D} = \mathcal{T}_{d-1} - \bar{\mathcal{H}}_d$ . If the difference has negative coefficients, the polynomial  $\mathcal{P}$  does not exist. Otherwise, since  $\mathcal{M}$  and  $\mathcal{D}$  are determined, the homogeneous component  $\mathcal{P}_d$  is uniquely determined, and we repeat the process.  $\square$

**Example 1 (continued).** Let  $\mathcal{T} = 2ab + a^2b + a^3b$ ,  $\mathcal{H} = a^3b + 2a^4b$ ,  $p = (23, 10)$  and  $\mu = a^2b$ . The polynomial  $\mathcal{P}$  has degree 2. Set  $\mathcal{P}_2 = a^2b\mathcal{D} + \mathcal{M}$ . Since  $\mathcal{H}$  has degree 1,  $\mathcal{M} = 0$ . Next,  $\mathcal{D} = \bar{\mathcal{P}}_2 = \mathcal{T}$ . Thus  $\mathcal{P}_2 = 2a^2bab + a^2ba^2b + a^2ba^3b$ . Set  $\mathcal{P}_1 = h * a^2b + \mathcal{Q}$ . From  $\hat{\mathcal{P}} = \mathcal{Q} + a^3b = \mathcal{H}$ , one gets  $\mathcal{Q} = 2a^4b$ . Content-Length: 12 821.

From  $\hat{\mathcal{P}} = \mathcal{Q} + a^3b = \mathcal{H}$ , one gets  $\mathcal{Q} = 2a^4b$ . Thus,  $\mathcal{P} = h * a^2b + 2a^2bab + a^2ba^2b + a^2ba^3b + 2a^4b$ , whence  $\text{par}(\mathcal{P}) = h * (2, 1) + (23, 10)$ . Thus  $h = 0$  and  $\mathcal{P} = 2a^2bab + a^2ba^2b + a^2ba^3b + 2a^4b$ .

We now consider the symmetric case, which is more involved. We will consider both polynomial and their reversals. More precisely, for any polynomial  $\mathcal{Q}$ , we denote  $\tilde{\mathcal{Q}}$  the polynomial  $\tilde{\mathcal{Q}} = \sum (\mathcal{Q}, x)\bar{x}$ . Observe that neither the shadow nor the trim reduction commutes with reversal. As an example, consider the polynomial

$$\mathcal{Q} = 2a + 2a^2baba + a^2ba^2ba^2 + 2a^4ba$$

of Example 2 above. We get

$$\tilde{\mathcal{Q}} = 2a + 2ababa^2 + a^2ba^2ba^2 + 2aba^4$$

from which  $\mu_{\underline{a}} = a^2b$ ,  $\mu_{\underline{b}} = ab$ . Also  $\bar{\mathcal{Q}} = 2aba + a^2ba^2 + a^4$ ,  $\tilde{\bar{\mathcal{Q}}} = 2aba^2 + a^2ba^2 + a^{10}$  and  $\hat{\mathcal{Q}} = a^6 + 2a^4ba$ ,  $\tilde{\hat{\mathcal{Q}}} = a^{14} + a^2ba^2ba^2$ .

**Proposition 5.4.** *Given symmetric normalized polynomials  $\mathcal{H}, \mathcal{H}_{\sim}$  and  $\mathcal{T}, \mathcal{T}_{\sim}$ , a vector  $p$ , and two blocks  $\mu, \mu_{\sim}$ , there is at most one symmetric normalized effectively computable polynomial  $\mathcal{P}$  such that  $\hat{\mathcal{P}} = \mathcal{H}$ ,  $\tilde{\hat{\mathcal{P}}} = \mathcal{H}_{\sim}$ ,  $\bar{\mathcal{P}} = \mathcal{T}$ ,  $\tilde{\bar{\mathcal{P}}} = \mathcal{T}_{\sim}$ ,  $\mu_{\mathcal{P}} = \mu$ ,  $\mu_{\tilde{\mathcal{P}}} = \mu_{\sim}$  and  $\text{par}(\mathcal{P}) = p$ .*

**Proof.** The proof is by induction as for Proposition 5.3. However, the initial step of the induction is different.

If  $\text{deg}(\mathcal{T}) = 0$ , then  $\text{deg}(\mathcal{P}) = 1$  and  $\mathcal{P} = h * a$ , where  $h$  is given by the vector  $p$ .

Assume  $\text{deg}(\mathcal{T}) = 1$ . Then  $\text{deg}(\mathcal{P}) = 2$ , and  $\mu$  has to be a regular block. Set

$$\mathcal{P} = h * a + \mu\mathcal{C} + \mathcal{N},$$

where the polynomial  $\mathcal{C}$  has degree 1, the polynomial  $\mathcal{N}$  is homogeneous of degree 2 and all initial blocks of words in  $\mathcal{N}$  are greater than  $\mu$ . One gets  $\hat{\mathcal{P}} = H * a + \mathcal{N}$ , where  $H * a$  is obtained from normalization of the shadow reduction of  $h * a + \mu\mathcal{C}$ . This determines  $\mathcal{N}$  to be  $\mathcal{H}_2$ . Similarly, write

$$\tilde{\hat{\mathcal{P}}} = h * a + \mu_{\sim}\mathcal{D} + \mathcal{M}.$$

By the same argument,  $\mathcal{M}$  is the homogeneous component of degree 2 of  $\mathcal{H}_{\sim}$ . Taking the reversal of the first relation, we get

$$\tilde{\mathcal{C}}\tilde{\mu} + \tilde{\mathcal{N}} = \mu_{\sim}\mathcal{D} + \mathcal{M}. \tag{3}$$

Set  $\mu = a^N u$ , and  $\mu_{\sim} = a^M v$  for nonempty words  $u, v$  over  $A - \{a\}$ . By the definition of  $\mu$  and  $\mu_{\sim}$ , every word of degree 2 in the support of  $\mathcal{P}$  is of the form  $a^n w a^m$  with  $n \geq N$ ,  $m \geq M$  and  $w$  a word over  $A - \{a\}$ . Set also  $\mathcal{C} = \sum_{m \geq M} c_m a^m$ ,  $\mathcal{D} = \sum_{n \geq N} d_n a^n$ , where the coefficients  $c_n$  and  $d_n$  are nonnegative integers. Then

$$\tilde{\mathcal{C}}\tilde{\mu} = c_M a^M \tilde{u} a^N + \sum_{m > M} c_m a^m \tilde{u} a^N, \quad \mu_{\sim}\mathcal{D} = d_N a^M v a^N + \sum_{n > N} d_n a^M v a^n,$$

whence

$$\tilde{\mathcal{N}} + c_M a^M \tilde{u} a^N + \sum_{m > M} c_m a^m \tilde{u} a^N = d_N a^M v a^N + \sum_{n > N} d_n a^M v a^n + \mathcal{M}.$$

This entirely determines the coefficients  $c_m$  for  $m > M$  and  $d_n$  for  $n > N$ . Moreover, if the polynomial  $\mathcal{P}$  exists, then  $c_M = d_N$ . If  $\tilde{u} \neq v$ , then  $c_M = d_N = 0$ .

Thus, assume  $\tilde{u} = v$ . It remains to compute the coefficients  $h$  and  $c_M$ . From  $\tilde{\hat{\mathcal{P}}} = \bar{\mathcal{T}}$ , one gets for the coefficient  $t$ :

$$t = h + \sum_{m \geq M} c_m * m \tag{4}$$

Finally, the  $a$ -component  $p_a$  of the Parikh vector  $p = \text{par}(\mathcal{P})$  has the expression

$$p_a = h + \sum_{m \geq M} c_m * (m + N) + \text{par}(\mathcal{N})_a. \tag{5}$$

These two equations allow to compute  $h$  and  $c_M$ .

If  $\text{deg}(\mathcal{T}) > 1$ , the homogeneous parts of degree at least 3 are obtained as in Proposition 5.3, and the homogeneous parts of degree at most 2 are obtained as above.

□

**Example 2 (continued).** Consider the following data:  $\mathcal{H} = 6a + 2a^4ba$ ,  $\mathcal{H}_\sim = 14a + a^2ba^2ba^2$ ,  $\mathcal{T} = 4a + 2aba + a^2ba^2$ ,  $\mathcal{T}_\sim = 10a + 2aba^2 + a^2ba^2$ ,  $\mu = a^2b$ ,  $\mu_\sim = ab$ ,  $p = (26, 8)$ . Since  $\mathcal{T}$  has degree 2, the polynomial  $\mathcal{P}$  we look for has degree 3. We set

$$\mathcal{P} = h * a + \mu\mathcal{C} + \mathcal{N} + \mu\mathcal{E} + \mathcal{R}$$

with  $\text{deg}(\mathcal{R}) = 3$ ,  $\text{deg}(\mathcal{E}) = \text{deg}(\mathcal{N}) = 2$ ,  $\text{deg}(\mathcal{C}) = 1$ . From  $\hat{\mathcal{P}} = \dots + \mathcal{R} = \mathcal{H}$ , we get  $\mathcal{R} = 0$ , and from  $\tilde{\mathcal{P}} = \dots + \mathcal{E} = \mathcal{T}$  we get  $\mathcal{E} = 2aba + a^2ba^2$ , thus the component  $\mathcal{P}_3 = \mu\mathcal{E} + \mathcal{R} = \mu\mathcal{E}$  is  $\mathcal{P}_3 = 2a^2baba + a^2ba^2ba^2$ . Consider again  $\hat{\mathcal{P}} = h * a + \nu(\mathcal{C}) + \mathcal{N} + 4a = \mathcal{H}$ . The coefficient  $4a$  comes from  $4a = \widehat{\mu\mathcal{E}}$ . Thus  $\mathcal{N} = 2a^4ba$ . Set  $\tilde{\mathcal{P}} = h * a + \mu_\sim \mathcal{D} + \mathcal{M} + \tilde{\mathcal{P}}_3$ . Then  $\hat{\tilde{\mathcal{P}}} = h * a + \nu(\mathcal{D}) + \mathcal{M} + 4a + a^2ba^2ba^2 = \mathcal{H}_\sim$  showing that  $\mathcal{M} = 0$ . It follows that  $\tilde{\mathcal{N}} + \tilde{\mathcal{C}}\tilde{\mu} = \mu_\sim \mathcal{D}$ . This gives the equation

$$2aba^4 + c_1aba^2 + \sum_{m > 1} c_m a^m ba^2 = d_2aba^2 + \sum_{n > 2} d_n aba^n$$

which implies  $\mathcal{C} = c_1a$ ,  $\mathcal{D} = d_2a^2 + 2a^4$  and  $\mathcal{P} = h * a + c_1a^2ba + 2a^4ba + 2a^2baba + a^2ba^2ba^2$ . The coefficient of  $a$  in  $\tilde{\mathcal{P}}$  is  $h + c_1 + 2$ , and since  $(\mathcal{P}, a) = (\mathcal{T}, a) = 4$ , we get  $h + c_1 = 2$ . Next, the  $a$ -coordinate of the Parikh vector  $\text{par}(\mathcal{P})$  is  $h + 3c_1 + 10 + 8 + 6$  and must be 26. This gives the equation  $2 = h + 3c_1$ . It follows that  $h = 2$  and  $c_1 = 0$ . Thus, our solution is precisely the polynomial  $\mathcal{Q}$  of the example.

### 5.3. Proof of the main theorem

We are now able to give the proof of Theorem 1.1.

Let  $\mathcal{P}$  and  $\mathcal{Q}$  be normalized polynomials such that  $\sqcup\sqcup(\mathcal{P}) = \sqcup\sqcup(\mathcal{Q})$ , and set  $S = \sqcup\sqcup(\mathcal{P}) = \sqcup\sqcup(\mathcal{Q})$ . In view of Corollary 3.6, we may assume that  $S$  is simple or symmetric. It follows that  $\mathcal{P}$  and  $\mathcal{Q}$  are simple or are symmetric. Also, it is clear that  $\text{par}(\mathcal{P}) = \text{par}(\mathcal{Q})$ .

Assume first that  $\mathcal{P}$  and  $\mathcal{Q}$  are  $a$ -simple. The proof is by induction on this Parikh vector. By Lemma 5.1, one gets  $\sqcup\sqcup(\tilde{\mathcal{P}}) = \sqcup\sqcup(\tilde{\mathcal{Q}})$ , and by Lemma 5.2, one gets  $\sqcup\sqcup(\tilde{\tilde{\mathcal{P}}}) = \sqcup\sqcup(\tilde{\tilde{\mathcal{Q}}})$ . As observed, the Parikh vectors of  $\tilde{\tilde{\mathcal{P}}}$  and of  $\tilde{\tilde{\mathcal{Q}}}$  are less than the Parikh vector of  $\mathcal{P}$ . By induction,  $\tilde{\tilde{\mathcal{P}}} = \tilde{\tilde{\mathcal{Q}}}$  and  $\tilde{\tilde{\mathcal{P}}} = \tilde{\tilde{\mathcal{Q}}}$ . By Proposition 5.3,  $\mathcal{P} = \mathcal{Q}$ .

The same argument holds if  $\mathcal{P}$  and  $\mathcal{Q}$  are symmetric, using Proposition 5.4. This completes the proof. □

## 6. A large example

Over the alphabet  $A = \{a, b\}$ , consider the following set  $S$  of 58 words given in alphabetical order (for the moment, forget about the marks  $*$  and  $+$ ):

```
+ aaaaababbb  aaaaabbabb  aaaaabbbab
+ aaaabaabbb  aaaabababb  aaaababbab  *aaaababbba  aaaabbaabb
  aaaabbabab  *aaaabbabba  aaaabbbbaab  *aaaabbbaba
+ aaabaaabbb  aaabaababb  aaabaabbab  *aaabaabbba  aaababaabb
  aaabababab  *aaabababba  aaababbaab  *aaababbaba  aaabbaaabb
  aaabbaabab  *aaabbaabba  aaabbabaab  *aaabbababa
+ aabaaaabbb  aabaaababb  aabaaabbab  *aabaaabbba  aabaaabaabb
  aabaababab  *aabaababba  aabaabbaab  *aabaabbaba  aababaaabb
  aababaabab  *aababaabba  aabababaab  *aababababa  aabbaaabab
  aabbaabaab  *aabbaababa
  abaaaababb  abaaaabbab  abaaabaabb  abaaababab  *abaaababba
  abaaabbaab  *abaaabbaba  abaabaaabb  abaabaabab  *abaabaabba
  abaababaab  *abaabababa  ababaaabab  ababaabaab  *ababaababa.
```

We start by some preliminary observations:

- All words of  $S$  contain 6 letters  $a$  and 4 letters  $b$ .
- All words of  $S$  start with the letter  $a$ .
- Some words end with a letter  $a$ , others with a  $b$ .

So, we do not need to split  $S$  according to the first letter; letter  $a$  will be the pivot. On the other hand, we have to split  $S$  into an  $a$ -simple set, say  $T = S[b]$  and an  $a$ -symmetric set, say  $U = S[a]$ . In view of the definition in Section 3, we look at the maximal number of  $a$ 's at the end of words in  $S$ . This number is 1, hence  $T = Sa^{-1}$ . This set is exactly the set of the 19 words marked with a  $*$  in  $S$ .

Similarly, we look at the maximal number of  $b$ 's at the end of words in  $S$ . This number is 3, hence  $U = Sb^{-3}$ . This set is exactly the set of four words marked with a  $+$  in  $S$ . In view of Theorem 3.5, we consider separately these two sets. If they are shuffle sets, then we get shuffle roots for them. These will be later composed to get the unique candidate for a shuffle root for  $S$ .

### 6.1. Computing the shuffle root of $U$

The set  $U$  is

*aaaaaba aaaabaa aaabaaa aabaaaa.*

We look for a polynomial  $\mathcal{U}$  such that  $U = \sqcup\sqcup(\mathcal{U})$ . For this, we apply shadow and trim reduction, as described in Section 5.1. The minimal initial block of  $U$  is  $\mu_U = aab$ . Since there is only one word in  $U$  beginning by  $\mu$ , the set  $\hat{U}$  is composed of the single word *aaaa*. Similarly, the leading block  $\lambda_U$  is  $a^5b$ ; it gives raise to  $\bar{U} = \{a\}$ . Hence, the polynomial  $\mathcal{U}$  associated to  $U$  will look like  $\mathcal{U} = aaba^i + ja$  with  $i + j = 4$

(= the number of  $a$  in words of  $U$  minus the number of  $a$  in  $\mu$ ). In order to compute  $i$  and  $j$ , we look at the reversal of  $U$ . The minimal initial block is  $\mu_{\tilde{U}} = ba$ . This yields  $\tilde{\tilde{U}} = \{aaaaa\}$ . The leading block  $\lambda_{\tilde{U}}$  of  $\tilde{U}$  is  $aaaab$ , giving raise to  $\hat{\tilde{U}} = \{aa\}$ . So, we get the new description of  $\mathcal{U}$  as  $a^kba + sa$  with  $k + s = 5$ . So, we get immediately

$$\mathcal{U} = aaba + 3a.$$

### 6.2. Computing the shuffle root of $T$

Here is the set  $T$  again:

*aaaababbb aaaabbabb aaaabbbab  
 aaabaabbb aaabababb aaababbab aaabbaabb aaabbabab  
 aabaaabbb aabaababb aabaabbab aababaabb aabababab aabbaabab  
 abaaababb abaaabbab abaabaabb abaababab ababaabab.*

Again, we look for a polynomial  $\mathcal{T}$  such that  $T = \sqcup\sqcup\mathcal{T}$ . In order to construct  $\tilde{\tilde{T}}$ , we look for the leading block in  $T$ ; it is  $\lambda_T = a^4b^3$ . The left quotient by this block yields  $\tilde{\tilde{T}} = \{ab\}$ . Hence,  $\tilde{\tilde{\mathcal{T}}} = ab$ .

In order to construct  $\hat{\hat{T}}$ , we look for the minimal initial block in  $T$ ; it is  $\mu = ab$ . The prefix used then to quotient  $T$  is  $abab$  yielding  $\hat{\hat{T}} = \{aabab\}$ . Hence,  $\hat{\hat{\mathcal{T}}} = aabab$ . Knowing  $\mu$ ,  $\tilde{\tilde{\mathcal{T}}}$  and  $\hat{\hat{\mathcal{T}}}$ , we construct  $\mathcal{T}$ . It has the form  $\mathcal{T} = aabab + h * ab$  for some integer  $h$ . Using the Parikh vector  $\text{par}(T) = (5, 4)$  of  $T$ , we get  $h = 2$ ; hence,

$$\mathcal{T} = aabab + 2ab.$$

### 6.3. Constructing the polynomial $\mathcal{S}$

Knowing  $\mathcal{U} = aaba + 3a$  and  $\mathcal{T} = aabab + 2ab$ , we may try now to construct  $\mathcal{S}$ . The longest word is  $aabab$ ; it is in  $\mathcal{S}$  and is responsible for  $aaba$  in  $\mathcal{U}$ . So, we are left with  $3a$  and  $2ab$ . This implies that  $2ab$  is in  $\mathcal{S}$  and leaves  $a$  alone in  $\mathcal{U}$ . Hence,

$$\mathcal{S} = aabab + 2ab + a.$$

The final step is now to compute the shuffle of  $\mathcal{S}$  and to check that it is precisely  $S$ . Note that this last step *must* be performed. Indeed, if you suppress from  $S$  any unmarked word, the resulting polynomial  $\mathcal{S}$  remains unchanged. Similarly, if you add to  $S$  a word such as  $aaabbbbaab$ , then

- it has the correct Parikh vector,
- it does not end by  $a$ , so it will not appear in  $T$ ,
- it does not end by  $bbb$ , so that it will not appear in  $U$ .

Hence, the computations of  $\mathcal{U}$ ,  $\mathcal{T}$  and  $\mathcal{S}$  will be the same as above.

## References

- [1] C. Allauzen, Calcul efficace du shuffle de  $k$  mots, Tech. Report IGM2000-02, Institut Gaspard Monge, Université Marne-la-Vallée, 2000.
- [2] M. Lothaire, Combinatorics on Words, Encyclopaedia of Mathematics and its Applications, vol. 17, Addison-Wesley, Reading, MA, 1983 (Reprinted in the Cambridge Mathematical Library, Cambridge University Press, Cambridge, 1997).
- [3] M. Nivat, G.D. Ramkumar, C. Pandu Rangan, A. Saoudi, R. Sundaram, Efficient parallel shuffle recognition, *Parallel Process. Lett.* 4 (1994) 455–463.
- [4] B. Pradeep, C. Murthy, A constant time string shuffle algorithm on reconfigurable meshes, *Internat. J. Comput. Math.* 68 (1998) 251–259.
- [5] D.E. Radford, A natural ring basis for the shuffle algebra and an application to group schemes, *J. Algebra* 58 (1979) 432–454.
- [6] C. Reutenauer, *Free Lie Algebras*, London Mathematical Society Monographs, Oxford Science Publications, Oxford, 1993.
- [7] J.-C. Spehner, Le calcul rapide des mélanges de deux mots. (Fast computing of the shuffle of two words), *Theoret. Comput. Sci.* 47 (1986) 181–203.