# On the Length of the Minimum Solution of Word Equations in One Variable

Kensuke Baba[1], Satoshi Tsuruta[2],
Ayumi Shinohara[1,2], and Masayuki Takeda[1,2]

[1] PRESTO, Japan Science and Technology Corporation (JST)
[2] Department of Informatics, Kyushu University 33, Fukuoka 812-8581, Japan
{baba,s-tsuru,ayumi,takeda}@i.kyushu-u.ac.jp

**Abstract.** We show the *tight upperbound* of the length of the minimum solution of a word equation $L = R$ in one variable, in terms of the differences between the positions of corresponding variable occurrences in $L$ and $R$. By introducing the notion of difference, the proof is obtained from Fine and Wilf's theorem. As a corollary, it implies that the length of the minimum solution is less than $N = |L| + |R|$.

## 1  Introduction

Word equations can be used to describe several features of strings, for example, they generalize pattern matching problem [3,4] with variables. Fig. 1 shows an example of word equations. The fundamental work in word equations is Makanin's algorithm [10] which decides whether a word equation has a solution (see for a survey on this topic [9]). Plandowski [11] introduced a PSPACE algorithm which gives the best upperbound so far known. On the other hand, the problem is known to be NP-hard [1]. An approach to the problem is to analyze word equations with a restricted number of variables. Charatonik and Pacholski [2], and Ilie and Plandowski [7] introduced a polynomial time algorithm for word equations in two variables. As to word equations in one variable, there is an efficient algorithm by Obono et al. [6] which solves a word equation $L = R$ in $O(N \log N)$ time in terms of $N = |L| + |R|$. Dąbrowski and Plandowski [5] presented an algorithm of $O(N + \sharp_x \log N)$ time complexity for the number $\sharp_x$ of occurrences of the variable $x$.

However, the upperbound of the length of the minimum solution of word equations is not exactly understood even for one-variable version. Let $\chi$ be the upperbound, that is, a word equation has a solution if and only if there exists a solution $A$ of length $|A| \leq \chi$. For any word equation in one variable, we can choose a single candidate for the solution of a length, therefore we have only to check for the $\chi$ candidates at most to decide whether a word equation has a solution. Indeed no $\chi$ leads a better result for the complexity as long as it is proportional to $N$, but from a practical viewpoint, $\chi$ is quite important. In [6], $\chi$ is taken to be equal to $4N$ without precise proof. Hence, we need to reduce $\chi$ as small as possible and prove it formally.

Let a, b be characters and $x$ be a variable. The word equation

$$xx\mathtt{baababa} = \mathtt{ababa}x\mathtt{ab}x$$

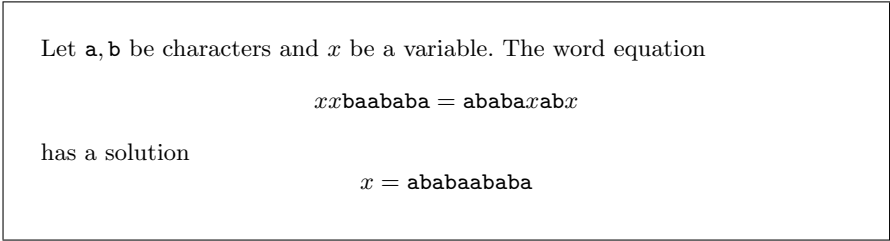has a solution

$$x = \mathtt{ababaababa}$$

**Fig. 1.** An example of word equation in one variable

In this paper, we show the *tight upperbound* of the minimum solution for one variable word equations, by introducing a new measure in terms of the positions of variable occurrences. The bound reveals that $\chi$ is less than $N$.

We now explain the basic idea briefly. A word equation in one variable is non-trivial only if both side of the equation have the same number of occurrences of the variable: Otherwise, the length of a possible solution is exactly determined by an integer equation on both the length of instance and the number of variable occurrences. Let $m$ be the number of occurrences. We focus on the fact that, for a word equation $L = R$, the "gap" between the $k$-th occurrence of the variable $x$ in $L$ and the $k$-th occurrence in $R$ is preserved for any substitution of a string $A$, as the gap between the corresponding occurrences of $A$ in $L[A/x]$ and $R[A/x]$. We denote the gaps by $d_k$ ($1 \le k \le m$). In the example in Fig. 1, $d_1 = 5$ and $d_2 = 7$. By utilizing this notion, the proof of the upperbound is essentially reducible to one for a word equation which has only one occurrence of $x$ in both side respectively. If $A$ is a solution and is longer than $d_k$, then the $k$-th pair of occurrences of $A$ overlap each other, that is, $d_k$ is a period of $A$. Therefore, by Fine and Wilf's theorem [9], the upperbound is $\max_{1 \le k \le m}\{d_k + p - \gcd(d_k, p)\} - 1$ for a period $p$ of $A$. Since the minimum length of $p$ is not larger than $\min_{1 \le k \le m, d_k \ne 0} d_k$, the tight upperbound will be given as $\max_{1 \le k \le m} d_k + \min_{1 \le k \le m, d_k \ne 0} d_k - 2$. Obviously, $\min_{1 \le k \le m, d_k \ne 0} d_k \le \max_{1 \le k \le m} d_k < |L|$. Thus $\chi$ is less than $N = 2|L|$.

## 2   Preliminaries

Let $\Sigma$ be an *alphabet* and $x \notin \Sigma$ be a *variable*. The empty word is denoted by $\varepsilon$. The length of a word $w$ is denoted by $|w|$, where $|\varepsilon| = 0$ and $|x| = 1$. The $i$-th element of a word $w$ is denoted by $w[i]$ for $1 \le i \le |w|$. The word $w[i]w[i+1] \cdots w[j]$ is called a *subword* of $w$, and denoted by $w[i : j]$. In particular, it is called a *prefix* if $i = 1$ and a *suffix* if $j = |w|$. For convenience, let $w[i : j] = \varepsilon$ for $j < i$.

A *period* of a non-empty word $w$ is defined as an integer $0 < p \le |w|$, such that $w[i] = w[i+p]$ for any $1 \le i \le |w| - p$. Note that the $|w|$ is always a period of $w$.

**Proposition 1 (Fine and Wilf).** *Let $p, q$ be periods of a word $w$. If $|w| \ge p + q - \gcd(p, q)$, then $\gcd(p, q)$ is also a period of $w$.*

A *word equation (in one variable)* is a pair of words over $\Sigma \cup \{x\}$ and is usually written by connecting two words with "=". A *solution* of a word equation $L = R$ is a homomorphism $\sigma : (\Sigma \cup \{x\})^* \to \Sigma^*$ leaving the letters of $\Sigma$ invariant and such that $\sigma(L) = \sigma(R)$. Since the solution is uniquely decided by a mapping of $x$ into $\Sigma^*$, in this paper we define a solution as a word $A \in \Sigma^*$ such that $A = \sigma(x)$. Therefore, we can rewrite the condition that $\sigma(L) = \sigma(R)$ by $L[A/x] = R[A/x]$, where the result $w[A/x]$ of the substitution of $A$ to $x$ in a word $w$ is defined inductively as:
if $w = \varepsilon$,      $w[A/x] = \varepsilon$;
if $w = a \in \Sigma$, $w[A/x] = a$;
if $w = x$,      $w[A/x] = A$;
if $w = w_1 w_2$, $w[A/x] = w_1[A/x]w_2[A/x]$.

If two words $L$ and $R$ have the same prefix $M$, the solution of a word equation $L = R$ is obtained by solving the word equation $L' = R'$ where $L = ML'$ and $R = MR'$. Therefore, we can assume without loss of generality that any word equation is of the form $xL_1 = BxR_1$ for a non-empty word $B$ which has no variable and words $L_1, R_1$. This form implies that any solution $A$ is a prefix of the word $B^k$ for a natural number $k$. By a similar argument for suffix, we can assume that either $L_1$ or $R_1$ ends with $x$. In particular, if $L$ and $R$ have exactly one occurrence of $x$ respectively, the word equation $L = R$ can be reduced to the form $xC = Bx$ for non-empty words $B, C$ which have no variable.

We denote by $\sharp_x(w)$ the number of occurrences of the variable $x$ in a word $w$. If a word equation $L = R$ has a solution $A$, the length of $L[A/x]$ is same as the length of $R[A/x]$. Hence we have $|L| + \sharp_x(L) \cdot (|A| - 1) = |R| + \sharp_x(R) \cdot (|A| - 1)$, and therefore $|A| = \frac{|L|-|R|}{\sharp_x(R)-\sharp_x(L)} + 1$. If $\sharp_x(L) \neq \sharp_x(R)$, the length of the solution is determined uniquely to the word equation and its upperbound is $|\,|L|-|R|\,| + 1 \leq \max(|L|, |R|)$. If $\sharp_x(L) = \sharp_x(R)$, we have $|L| = |R|$.

**Proposition 2 ([6]).** *Let $L = R$ be a word equation.*
*(i) If $\sharp_x(L) \neq \sharp_x(R)$, the length of the solution is determined uniquely with respect to $L = R$ and is at most $\max(|L|, |R|)$.*
*(ii) If $\sharp_x(L) = \sharp_x(R)$, $L = R$ has a solution only if $|L| = |R|$.*

## 3   Solutions

We show the upperbound of the length of the minimum solution of word equations in one variable. By Proposition 2, we have only to consider the word equation $L = R$ in the situation that $\sharp_x(L) = \sharp_x(R)$ and $|L| = |R|$. Let $m = \sharp_x(L) = \sharp_x(R)$ and $n = |L| = |R|$. We denote by $\ell_1^x, \cdots, \ell_m^x$ and $r_1^x, \cdots, r_m^x$ the positions of occurrences of $x$ in $L$ and $R$, respectively in increasing order. We define $\ell_k^A$ and $r_k^A$ for a word $A$ and $1 \leq k \leq m$ as

$$\ell_k^A = \ell_k^x + (k-1)(|A| - 1),$$
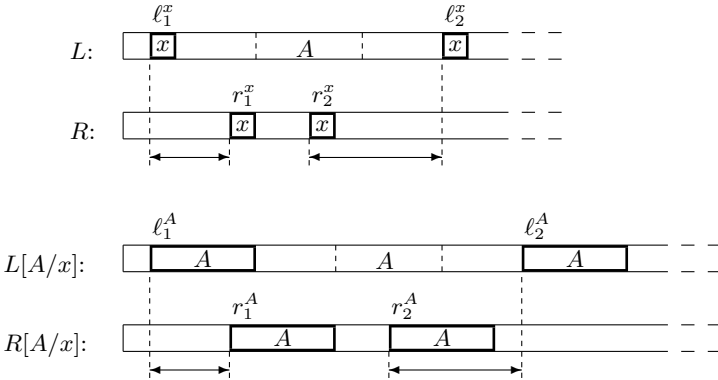$$r_k^A = r_k^x + (k-1)(|A| - 1).$$

**Fig. 2.** The difference $\ell_k^x - r_k^x$ is equal to the difference $\ell_k^A - r_k^A$ for any $A$

$\ell_k^A$ is, intuitively, the position in $L[A/x]$ of a occurrence of $A$ substituted to the $k$-th occurrence of $x$ in $L$ (which is not always the $k$-th occurrence). Therefore, $\ell_k^A - r_k^A$ is the difference between it and the position of the corresponding occurrence of $A$ in $R[A/x]$. The difference does not depend on the length of $A$, see Fig. 2.

**Proposition 3.** *For any word $A$, any word equation $L = R$, and integer $1 \le k \le m$,*
(i) $\ell_k^A - r_k^A = \ell_k^x - r_k^x$,
(ii) $L[A/x][\ell_k^A : \ell_k^A + |A| - 1] = R[A/x][r_k^A : r_k^A + |A| - 1] = A$.

*Proof.* (i) Trivial by the definition.
(ii) We prove for $L$. By the definition of substitution, $L[A/x]$ is represented as

$$L[A/x] = L[1 : \ell_1^x - 1]AL[\ell_1^x + 1 : \ell_2^x - 1] \cdots$$
$$L[\ell_{k-1}^x + 1 : \ell_k^x - 1]AL[\ell_k^x + 1 : \ell_{k+1}^x - 1] \cdots$$
$$L[\ell_{m-1}^x + 1 : \ell_m^x - 1]AL[\ell_m^x + 1 : n]. \tag{1}$$

The length of the prefix of $L[A/x]$ which ends $L[\ell_{k-1}^x + 1 : \ell_k^x - 1]$ equals to $(\ell_1^x - 1) + \sum_{i=2}^{k}\{(\ell_i^x - 1) - (\ell_{i-1}^x + 1) + 1\} + (k-1)|A| = \ell_k^x - k + (k-1)|A| = \ell_k^A - 1$ for any $1 \le k \le m$. Thus $\ell_k^A$ is the position of the occurrence of $A$ which is the next to $L[\ell_{k-1}^x + 1 : \ell_k^x - 1]$ in the right side of Eq. (1). □

We denote by $d_k$ the absolute value of the difference, that is,

$$d_k = |\ell_k^x - r_k^x|$$

for $1 \le k \le m$. Then we have the following lemma.

**Lemma 1.** *Let $A$ be a solution of a word equation $L = R$. For $1 \le k \le m$ and $d_k \ne 0$, if $|A| \ge d_k$ then $A$ has a period $d_k$.*

*Proof.* We can assume $r_k^x < \ell_k^x$ without loss of generality. If $|A| = d_k$, by the definition, $d_k$ is a period of $A$. If $|A| > d_k$, by Proposition 3 (i), $\ell_k^A = r_k^A + \ell_k^x - r_k^x = r_k^A + d_k < r_k^A + |A|$. Since $A$ is a solution of $L = R$, we consider subwords of $L[A/x]$ and $R[A/x]$, then $L[A/x][\ell_k^A : r_k^A + |A| - 1] = R[A/x][\ell_k^A : r_k^A + |A| - 1]$. By Proposition 3 (ii), $L[A/x][\ell_k^A : r_k^A + |A| - 1] = A[1 : |A| - (\ell_k^A - r_k^A)]$ and $R[A/x][\ell_k^A : r_k^A + |A| - 1] = A[1 + (\ell_k^A - r_k^A) : |A|]$. Thus, by Proposition 3 (i), $A[1 : |A| - (\ell_k^x - r_k^x)] = A[1 + (\ell_k^x - r_k^x) : |A|]$ which implies that $\ell_k^x - r_k^x$ is a period of $A$. □

**Lemma 2.** *Let $A$ be a solution of a word equation $L = R$ and $p$ be a period of $A$. If*

$$|A| \geq \max_{1 \leq k \leq m} d_k + p - 1,$$

*then the prefix $A[1 : |A| - p]$ of $A$ is also a solution of $L = R$.*

*Proof.* We prove by induction on the number $m = \sharp_x(L) = \sharp_x(R)$.
(Base step) By the argument in Section 2, we can assume $L = xC$ and $R = Bx$ with $B, C \in \Sigma^+$. By Lemma 1, $d_1 = |B|$ is a period of $A$. By Proposition 1, $\gcd(d_1, p)$ is a period of $A$, moreover it is also a period of $AC$ and $BA$. Since $A[1 + \gcd(d_1, p) : |A|] = A[1 : |A| - \gcd(d_1, p)]$, we have

$$\begin{aligned}
A[1 : |A| - k\gcd(d_1, p)]C &= (AC)[1 + k\gcd(d_1, p) : |A|] \\
&= (BA)[1 + k\gcd(d_1, p) : |A|] \\
&= BA[1 : |A| - k\gcd(d_1, p)]
\end{aligned}$$

for a natural number $k$ such that $k\gcd(d_1, p) \leq |A|$.
(Induction step) We can assume $L = L'xC$ and $R = R'xBx$ with $L', R' \in (\Sigma \cup \{x\})^+$ and $B, C \in \Sigma^+$. Then we have $d_m = |C|$ and $L'[A/x]AC = R'[A/x]ABA$. If $|C| \leq |B|$, the result is obviously obtained by induction for two equations $L' = R'xB[1 : |B| - |C|]$ and $xC = B[|B| - |C| + 1 : |B|]x$. If $|C| > |B|$, we have $|ABA| > |AC| > |BA|$ by the assumption $|A| \geq \max_{1 \leq k \leq m} d_k + p - 1$. Hence the occurrence of $A$ starting at $\ell_m^A$ in $L[A/x]$ and the occurrence of $A$ starting at $r_{m-1}^A$ in $R[A/x]$ have a non-trivial overlapping $Q$. (This situation is illustrated in Fig. 3.) Now we consider two equations $L'Q = R'x$ and $xC = QBx$. The assumption $L'[A/x]AC = R'[A/x]ABA$ implies $L'[A/x]Q = R'[A/x]A$ and $AC = QBA$, that is, $A$ is a solution of the equations. Then, by induction hypothesis, we have $L'[A'/x]Q = R'[A'/x]A'$ and $A'C = QBA'$ where $A' = A[1 : |A| - p]$. Thus, we have $L[A'/x] = L'[A'/x]A'C = L'[A'/x]QBA' = R'[A'/x]A'BA' = R[A'/x]$. □

**Theorem 1 (Tight upperbound).** *For any word equation $L = R$ such that $\sharp_x(L) = \sharp_x(R)$, the length of the minimum solution is at most*

$$\max_{1 \leq k \leq m} d_k + \min_{1 \leq k \leq m, d_k \neq 0} d_k - 2.$$
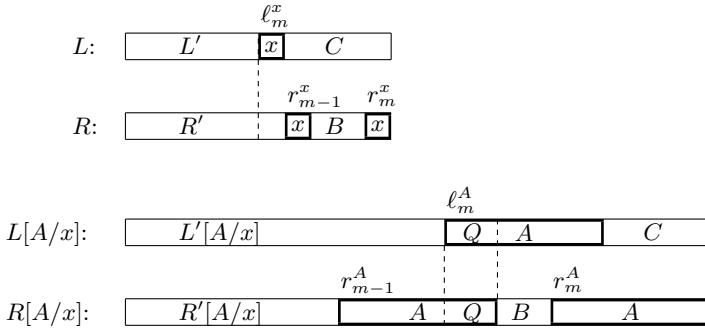
*The bound is tight.*

**Fig. 3.** If $|C| > |B|$ and $|A| \geq d_m = |C|$, then $|ABA| > |AC| > |BA|$ and two occurrences of $A$ starting at $\ell_m^A$ and $r_{m-1}^A$ have a overlap $Q$

*Proof.* Assume a word equation has a solution $A$ such that $|A| \geq \max_{1 \leq k \leq m} d_k + \min_{1 \leq k \leq m, d_k \neq 0} d_k - 1$. By Lemma 1, $A$ has a period $p \leq \min_{1 \leq k \leq m, d_k \neq 0} d_k$. Hence, by Lemma 2, $A[1 : |A| - p]$ is also a solution of the word equation. Therefore $A$ is not the minimum solution.

To see that the bound is tight, let us consider the following word equation:

$$xx\mathtt{baababa} = \mathtt{ababa}x\mathtt{ab}x.$$

We can verify that the solution of length 10

$$x = \mathtt{ababaababa}.$$

is in fact the minimum solution. Since $d_1 = 5$ and $d_2 = 7$, we have $\max_{1 \leq k \leq 2} d_k = 7$ and $\min_{1 \leq k \leq 2} d_k = 5$. Thus $\max_{1 \leq k \leq 2} d_k + \min_{1 \leq k \leq 2} d_k - 2 = 10$, which shows the bound is tight. □

In case of binary alphabet, the minimum solution which length is the upper bound is *central* which is defined as:
*A word is central if and only if it is in the set*

$$0^* \cup 1^* \cup (P \cap P10P)$$

*where $P$ is the set of palindrome words.*
It is obtained by the proof of Lemma 2 and the fact that: a word $w$ is central if and only if it has two periods $p$ and $q$ such that $\gcd(p, q) = 1$ and $|w| = p + q - 2$ [9, pp. 69–70].

We also have the following relaxed upperbound, since $\min_{1 \leq k \leq m, d_k \neq 0} d_k \leq \max_{1 \leq k \leq m} d_k < |L|$.

**Corollary 1.** *For any word equation $L = R$ such that $\sharp_x(L) = \sharp_x(R)$, the length of the minimum solution is at most $N - 4 = |L| + |R| - 4$.*

Consequently, we have the following upperbound by Proposition 2.

**Corollary 2.** *For any word equation $L = R$, the length of the minimum solution is at most $N - 1$.*

**Table 1.** The numbers of *solvable* word equations in one variable in $\mathcal{E}$, classified by the lengths of their minimum solutions

| | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 |
|---|---|---|---|---|---|---|---|---|---|
| | | | | length of $L$ (and $R$) | | | | | |
| 0 | 4 | 32 | 220 | 1388 | 8364 | 49120 | 284204 | 1630124 | 9303292 |
| 1 | 4 | 20 | 104 | 548 | 2868 | 14856 | 76236 | 388212 | 1964612 |
| 2 | 0 | 12 | 56 | 252 | 1208 | 5844 | 28268 | 136536 | 657868 |
| 3 | 0 | 0 | 24 | 140 | 564 | 2488 | 11304 | 53008 | 250296 |
| 4 | 0 | 0 | 0 | 60 | 260 | 1148 | 4764 | 20784 | 95868 |
| 5 | 0 | 0 | 0 | 0 | 116 | 580 | 2052 | 8592 | 36076 |
| 6 | 0 | 0 | 0 | 0 | 8 | 264 | 1152 | 4368 | 16152 |
| 7 | 0 | 0 | 0 | 0 | 0 | 8 | 504 | 2148 | 7532 |
| 8 | 0 | 0 | 0 | 0 | 0 | 0 | 24 | 1084 | 4404 |
| 9 | 0 | 0 | 0 | 0 | 0 | 0 | 8 | 48 | 2120 |
| 10 | 0 | 0 | 0 | 0 | 0 | 0 | 8 | 36 | 136 |
| 11 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 8 | 24 |
| 12 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 13 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 8 |
| 14 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 8 |

## 4   String Statistics Problem

We are developing a system whose aim is to experimentally analyze the combinatorial property and structures of word equations. As a first step, we are recording *all* solvable word equations *(up to a moderate length)* in one variable together with their minimum solutions. By the fact that: for any word $w$, there exists a binary word $w'$ which has the same set of periods as $w$ [9, pp. 275–279], we have only to consider a binary alphabet to find out the relation between the length of an equation and the length of its solutions. For a fixed alphabet $\Sigma = \{\mathtt{a}, \mathtt{b}\}$ and a specified length $n$, we enumerate the set $\mathcal{E}$ of *all* word equations $L = R$ such that (1) both $\mathtt{a}$ and $\mathtt{b}$ appear either $L$ or $R$, (2) $|L| = |R| = n$, (3) $L$ and $R$ contains the same number of variables, and (4) the pairs $(L[1], R[1])$ and $(L[n], R[n])$ must be taken from $\{(x, \mathtt{a}), (x, \mathtt{b}), (\mathtt{a}, x), (\mathtt{b}, x)\}$.

Then for each word equation in $\mathcal{E}$, we try to find the minimum solution by checking each prefix of $B^k$ (where $B$ is a constant prefix of either $L$ or $R$) in increasing order up to $2n - 4$. If a solution is found, we logged it and turn to the next equation. Otherwise, we can conclude that the word equation has no solution, thanks to the upperbound we have shown (Corollary 1).

For interested readers, Table 1 shows the numbers of the solvable word equations in $\mathcal{E}$, classified by the lengths of their minimum solutions. At $i$-th row and column labeled $n = |L|$ of the table $T$, we fill the number of word equations in $\mathcal{E}$ of length $|L| = |R| = n$ whose minimum solution is of length $i$. Remark that some equations may be equivalent each other, by either replacing $\mathtt{a}$ with $\mathtt{b}$, exchanging left-side with right-side, or reversing the formulae. We did not exclude these duplications. For example, $T(0, 3) = 4$ corresponds to the number of equations $\{\mathtt{ab}x = x\mathtt{ab}, \mathtt{ba}x = x\mathtt{ba}, x\mathtt{ab} = \mathtt{ab}x, x\mathtt{ba} = \mathtt{ba}x\}$, where the

empty string is a solution to them. They are equivalent each other. Moreover, $T(1,3) = 4$ corresponds to $\{\mathtt{ab}x = x\mathtt{ba}, \mathtt{ba}x = x\mathtt{ab}, x\mathtt{ab} = \mathtt{ba}x, x\mathtt{ba} = \mathtt{ab}x\}$, whose minimum solutions are of length 1. They are essentially the same.

Let us pick up some *interesting* pairs of equation and its minimum solution.

- $\langle xx\mathtt{baababa} = \mathtt{ababa}x\mathtt{ab}x, \mathtt{ababaababa} \rangle$, from $T(10,9) = 8$, which was used to prove the tightness of the upperbound. This is a *unique* instance in $T(10,9) = 8$, since the other 7 instances are all equivalent to it.
- $\langle xx\mathtt{baabababa} = \mathtt{ababab}ax\mathtt{ab}x, \mathtt{abababaabababa} \rangle$, from $T(14,11) = 8$, which also matches the upperbound. This is a *unique* instance in $T(14,11) = 8$, since the other 7 instances are all equivalent to it.
- $\langle x\mathtt{ab}x\mathtt{baaaaaa} = \mathtt{aaaaaaba}x\mathtt{b}x, \mathtt{aaaaaabaaaaaa} \rangle)$. This is a *unique* instance in $T(13,11) = 8$.

## 5    Conclusion

We showed the *tight upperbound* of the length of minimum solution of word equations in one variable. The upperbound is easily computed from a given word equation. Moreover, we showed concrete examples which match the bound. As a corollary, we also have a more relaxed upperbound which is easier applicable: the length of the minimum solution is less than the size of the total length of a word equation.

Khmelevskiĭ [8, pp. 12] proved that if a word equation $C_0 x C_1 \cdots x C_u = x B_1 \cdots x B_v$ is solvable, it has a solution of length smaller than $M^2 + 3M$ where $M = \max_{i,j}\{u, v, |C_i|, |B_j|\}$. When we consider the upperbound in terms of the length $N$ of a given word equation, the order of this value comes up to $N^2$ since $M \leq N - 1$. Even for the original expression, we can show that the value $M^2 + 3M - 1$ never be less than the upperbound of our result for a non-trivial word equation. Let $\nu = u = v$ and $\lambda = \max_{i,j}\{|C_i|, |B_j|\}$. Then $M = \max\{\nu, \lambda\}$. By the definition of $d_k$, we have $\min_{k,d_k \neq 0} d_k \leq |C_0| \leq \lambda$ and $\max_k d_k \leq \max\{\sum_{i=0}^{k-1} |C_i|, \sum_{i=k}^{\nu} |C_i|\} \leq \nu\lambda$. Therefore, $\max_k d_k + \min_{k,d_k \neq 0} d_k - 2 \leq \nu\lambda + \lambda - 2 \leq M^2 + 2M - 2 \leq M^2 + 3M - 1$.

Thanks to the bound, we could perform a comprehensive analysis of word equations in one variable up to a moderate size the equations, by enumerating all word equations and solving them one by one. We showed some statistics of the lengths of minimum solutions.

## Acknowledgements

# References

1. Angluin, D.: Finding Patterns Common to a Set of Strings. J. Comput. Sys. Sci., Vol. 21 (1980) 46–62
2. Charatonik, W. and Pacholski, L.: Word Equations in Two Variables. Proc. IWW-ERT'91, LNCS, Vol. 677 (1991) 43–57
3. Crochemore, M. and Rytter, W.: Text Algorithms. Oxford University Press, New York (1994)
4. Crochemore, M. and Rytter, W.: Jewels of Stringology. World Scientific (2003)
5. Dąbrowski, R. and Plandowski, W.: On Word Equations in One Variable. Proc. MFCS2002, LNCS Vol. 2420 (2002) 212–220
6. Eyono Obono, S., Goralcik, P., and Maksimenko, M.: Efficient Solving of the Word Equations in One Variable. Proc. MFCS'94, LNCS Vol. 841 (1994) 336–341
7. Ilie, L. and Plandowski, W.: Two-Variable Word Equations. Proc. STACS2000, LNCS Vol. 1770 (2000) 122–132
8. Khmelevskiĭ, Yu.I.: Equations in Free Semigroups. Proc. Steklov Inst. of Mathematics 107, AMS (1976)
9. Lothaire, M.: Algebraic Combinatorics on Words. Cambridge University Press (2002)
10. Makanin, G.S.: The Problem of Solvability of Equations in a Free Semigroup. *Mat. Sb.* Vol. 103, No. 2, 147–236. In Russian; English translation in: *Math. USSR Sbornik*, Vol. 32 (1977) 129–198
11. Plandowski, W.: Satisfiability of Word Equations with Constants is in PSPACE. Proc. FOCS'99, IEEE Computer Society Press (1999) 495–500