

# SETS REPRESENTED AS THE LENGTH- $n$ FACTORS OF A WORD

Shuo Tan and Jeffrey Shallit  
School of Computer Science,  
University of Waterloo,  
Waterloo, Ontario  
N2L 3G1 Canada

s22tan@uwaterloo.ca shallit@cs.uwaterloo.ca

September 14, 2013

- Let  $F_n(w)$  denote the set of length- $n$  factors of a word  $w$ .
- Let  $CF_n(w)$  denote the set of length- $n$  cyclic factors of a word  $w$ .

## DEFINITION

A word  $w$  **witnesses** (resp., **cyclically witnesses**) a subset  $S$  of  $\Sigma^n$  if  $F_n(w) = S$  (resp.,  $CF_n(w) = S$ ).

A subset  $S$  of  $\Sigma^n$  is **representable** (resp. **cyclically representable**) if there exists a word  $w$  that witnesses (cyclically witnesses)  $S$ .

## EXAMPLE (CYCLIC REPRESENTABLE SET)

The word  $w = 000011$  cyclically witnesses the set  $\{000, 001, 011, 110, 100\}$ . Thus the set  $\{000, 001, 011, 110, 100\}$  is cyclically representable.

## EXAMPLE (REPRESENTABLE SET)

$w = 00001$ .  $F_3(w) = \{000, 001\}$ . Thus, the set  $\{000, 001\}$  is representable. The set  $\{000, 001\}$  is not cyclically representable.

## EXAMPLE (UNREPRESENTABLE SET)

The set  $\{00, 11\}$  is not representable.

## DEFINITION

A binary **De Bruijn sequence**  $B_n$  of order  $n$ , is a word  $w$  over  $\Sigma$  such that every possible word of length  $n$  appears exactly once as a cyclic factor of  $w$ .

## EXAMPLE (DE BRUIJN SEQUENCES)

$n$	$B_n$
2	0011
3	00010111
4	0000100110101111
5	00000111011010111110011000101001

## THEOREM

*For any order  $n > 0$ , a binary De Bruijn sequence  $B_n$  exists.*

## COROLLARY

*The set  $\Sigma^n$  is cyclically representable.*

## PROBLEM (1)

*How many subsets of  $\Sigma^n$  are (cyclically) representable?*

## PROBLEM (2)

*For how many subsets  $S$  of  $\Sigma^n$  does there exist a word  $w$  of length  $t$  that witnesses  $S$ ?*

# NUMERICAL RESULTS FOR PROBLEM 1

$n$	non-cyclic case	cyclic case
1	3	3
2	14	6
3	121	27
4	5921	972
5	20020315	2466131

# BOUNDS FOR PROBLEM 1

## THEOREM

*A lower bound on the number of cyclically representable subsets is  $2^{2^{n-1}}$ .*

## THEOREM

*An upper bound on the number of cyclically representable subsets is  $10^{2^{n-2}}$ .*

The number of subsets of  $\Sigma^n$  is  $2^{2^n} = 16^{2^{n-2}}$ .



## PROBLEM 2

Let  $T(n, t)$  denote the number of subsets  $S$  of  $\Sigma^n$  such that there exists a word  $w$  of length  $t$  that witnesses  $S$ .

PROBLEM (2)

$T(n, t) = ?$

For this problem, we have some results on the non-cyclic case.

# EXAMPLE

## EXAMPLE

We fix  $n = 2$  and  $t = 3$ .

word of length 3	witnesses
000	{00}
001	{00, 01}
010	{01, 10}
011	{01, 11}
100	{10, 00}
101	{01, 10}
110	{11, 10}
111	{11}

Note that 010 and 101 'coincide'. Thus  $T(2, 3) = 7$ .

## PROBLEM 2

In order to compute  $T(n, t)$ , we compute the number of coincidences; namely, we consider the number of words that witness the same subset of  $\Sigma^n$ .

Suppose  $S \subseteq \Sigma^n$ . Let  $C_t(S)$  denote the number of words of length  $t$  that witness  $S$ . Then we have

$$T(n, t) = 2^t - \sum_{S \in \Sigma^n, C_t(S) > 1} (C_t(S) - 1)$$

.

# SOME EASY CASES

1.  $T(n, n) = 2^n$  for  $n > 0$ .
2.  $T(n, n + 1) = 2^{n+1} - 1$  for  $n \geq 1$ .
3.  $T(n, n + 2) = 2^{n+2} - 5$  for  $n \geq 2$ .
4.  $T(n, n + 3) = 2^{n+3} - 14$  for  $n \geq 3$ .

## DEFINITION

A positive integer  $p$  is a **period** of a (finite) word  $w$ , if for any  $1 \leq i \leq |w| - p$  we have  $w[i] = w[i + p]$ .

Let  $\pi(w)$  denote the minimal period of  $w$ .

## EXAMPLE

Both 3 and 5 are periods of the word  $w = 010010$ . The minimal period of  $w$  is 3.

## DEFINITION

A **root**  $r(w)$  of a word  $w$  is the prefix of  $w$  with length  $\pi(w)$ . We say that two words  $w$  and  $w'$  are **root-conjugate** if  $r(w)$  and  $r(w')$  are conjugate.

## EXAMPLE

The words  $w_1 = 01001$  and  $w_2 = 10010$  are root-conjugate since  $r(w_1) = 010$  and  $r(w_2) = 100$  are conjugate.

## THEOREM

*Let  $t, n, k$  be integers such that  $t = n + k$ ,  $n \geq k + 1$ , and  $k \geq 0$ . For any distinct words  $w, w'$  of length  $t$ , we have  $F_n(w) = F_n(w')$  if and only if  $w$  and  $w'$  are root-conjugate and  $\pi(w) = \pi(w') \leq k + 1$ .*

## EXAMPLE

For example, let  $t = 7$ ,  $n = 4$ , and  $k = 3$ . We consider the words  $w_1 = 0110110$  and  $w_2 = 1011011$  that are root-conjugate. We have  $F_4(w_1) = F_4(w_2) = \{0110, 1101, 1011\}$ .

## THEOREM

*Let  $t, n, k$  be integers such that  $t = n + k$ ,  $n \geq k + 1$ , and  $k \geq 0$ . For any distinct words  $w, w'$  of length  $n$ , we have  $F_n(w) = F_n(w')$  if and only if  $w$  and  $w'$  are root-conjugate and  $\pi(w) = \pi(w') \leq k + 1$ .*

## PROOF.

If  $w$  and  $w'$  are root-conjugate with period  $p \leq k + 1$ , then there are  $p$  places to begin, and considering consecutive factors of length  $n + p - 1$  gives exactly  $p$  distinct length- $n$  factors.



## PROOF (CONT.)

For the other direction, we give a proof by induction on  $k$ .

- The base case where  $k = 0$  is easy. In this case  $t = n$  and thus  $F_n(w) = \{w\}$  and  $F_n(w') = \{w'\}$ . Thus  $w = w'$ , contradicting the fact that  $w$  and  $w'$  are distinct.
- We assume the result holds for  $k - 1$  and we prove it for  $k$ .
- Let  $p_i(w)$  denote the length- $i$  prefix of  $w$ . Let  $s_i(w)$  denote the length- $i$  suffix of  $w$ .



## PROOF (CONT.)

We consider the following cases:

1.
  - $p_{n-1}(w)$  appears only once as a factor of  $w$ .
  - In this case we can prove that  $p_n(w) = p_n(w')$ .
  - Let  $s = w[2..t]$  and  $s' = w[2..t]$ . We can show that  $F_n(s) = F_n(s')$ .
  - Let  $t' = t - 1$  and  $k' = k - 1$ . Applying induction we obtain that  $s$  and  $s'$  are root-conjugate and  $\pi(w) = \pi(w') \leq k' + 1 = k$ .
  - $w$  and  $w'$  are root-conjugate and  $\pi(w) = \pi(w') \leq k + 1$ .



2. Similarly we can handle the case where  $s_{n-1}(w)$  appears only once as a factor of  $w$ .
3. The last case is where both  $p_{n-1}(w)$  and  $s_{n-1}(w)$  appear more than once. We again consider two sub-cases depending on whether  $p_{n-1}(w) = s_{n-1}(w)$ .

## COROLLARY

For  $t \leq 2n - 1$ , we have  $T(n, t) = 2^t - \sum_{k=1}^{t-n+1} \frac{k-1}{k} \sum_{d|k} \mu\left(\frac{k}{d}\right) 2^d$ ,  
 where  $\mu(\cdot)$  is the Möbius function.

1.  $C_t(S) > 1$  iff there exists a word  $w$  that witnesses  $S$  with  $\pi(w) \leq k + 1$ .
2.  $C_t(S) = \pi(w)$ . Correspond to a set of root-cojugate words that witnesses  $S$ . Represented by their lexicographically least roots (Lyndon words).

As an example, let  $t = 5$  and  $n = 3$ . The set  $\{010, 100, 001\}$  is witnessed by 01001, 00100, 10010. It can be represented by 001.

Thus we have

$$\begin{aligned}
 T(n, t) &= 2^t - \sum_{S \in \Sigma^n, C_t(S) > 1} (C_t(S) - 1) \\
 &= 2^t - \sum_{\substack{w \text{ is a Lyndon word and} \\ \pi(w) \leq t - n + 1}} (\pi(w) - 1) \\
 &= 2^t - \sum_{i=1}^{t-n+1} (i-1) \cdot L(i)
 \end{aligned}$$

where  $L(i) = \frac{1}{i} \sum_{d|i} \mu(\frac{i}{d}) 2^d$  is the number of Lyndon words of length  $i$ .

- 1 Let  $\mathring{R}_n$  denote the set of all non-empty circularly representable subsets of  $\Sigma^n$ . Does the limit  $\lim_{n \rightarrow \infty} |\mathring{R}_n|^{\frac{1}{2n}}$  exist?
- 2 Derive a formula for  $T(n, t)$  where  $t = 2n$ .

We conjecture that if  $x$  and  $y$  are distinct binary words of length  $2n$  with  $F_n(x) = F_n(y)$  then  $\pi(x) = \pi(y)$  and furthermore  $x$  and  $y$  are root-conjugate. However, it is possible in this case that  $\pi(x) > n + 1$ .

Furthermore it seems that if  $\pi(x) > n + 1$ , then  $x = uv01vu$  and  $y = uv10v^R u$  (or vice versa) for some nonempty words  $u, v$  where  $u$  is the longest palindrome prefix of  $uv$  and  $\pi(x) = t - |v|$ .