

# Tighter Bounds for the Sum of Irreducible LCP Values

Juha Kärkkäinen<sup>1</sup>   Dominik Kempa<sup>1</sup>   Marcin Piątkowski<sup>2</sup>

<sup>1</sup> University of Helsinki



<sup>2</sup> Nicolaus Copernicus University



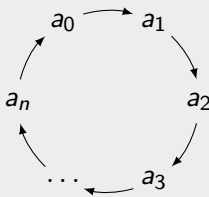
- 1 Cyclic words
- 2 Irreducible LCP values
- 3 Upper bound for the sum of irreducible values
- 4 Lower bound for the sum of irreducible values

$$W = \{v_1, v_2, v_3\} = \{aab, aba, baa\}$$

suf( $W$ )

---

$\langle 1, 0 \rangle$	<b>a a b</b> a a b $\dots$
$\langle 1, 1 \rangle$	<b>a b a</b> a b a $\dots$
$\langle 1, 2 \rangle$	<b>b a a</b> b a a $\dots$
$\langle 2, 0 \rangle$	<b>a a b</b> a a b $\dots$
$\langle 2, 1 \rangle$	<b>a b a</b> a b a $\dots$
$\langle 2, 2 \rangle$	<b>b a a</b> b a a $\dots$
$\langle 3, 0 \rangle$	<b>a b a</b> b a b $\dots$
$\langle 3, 1 \rangle$	<b>b a b a</b> b a $\dots$

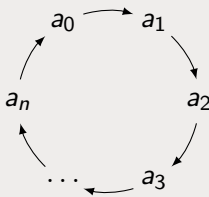


# Cyclic suffix array

$$W = \{v_1, v_2, v_3\} = \{aab, aab, ab\}$$

$SA_W$	$\text{suf}(W)$
--------	-----------------

$\langle 1, 0 \rangle$	<b>a a b</b> a a b ...
$\langle 2, 0 \rangle$	<b>a a b</b> a a b ...
$\langle 1, 1 \rangle$	<b>a b a</b> a b a ...
$\langle 2, 1 \rangle$	<b>a b a</b> a b a ...
$\langle 3, 0 \rangle$	<b>a b a</b> b a b ...
$\langle 1, 2 \rangle$	<b>b a a</b> b a a ...
$\langle 2, 2 \rangle$	<b>b a a</b> b a a ...
$\langle 3, 1 \rangle$	<b>b a b</b> a b a ...

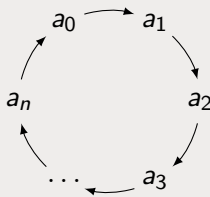


# Cyclic suffix array

$$W = \{v_1, v_2, v_3\} = \{aab, aab, ab\}$$

$SA_W$	$\text{suf}(W)$
--------	-----------------

$\langle 1, 0 \rangle$	<b>a a b</b> a a b $\dots$
$\langle 2, 0 \rangle$	<b>a a b</b> a a b $\dots$
$\langle 1, 1 \rangle$	<b>a b a</b> a b a $\dots$
$\langle 2, 1 \rangle$	<b>a b a</b> a b a $\dots$
$\langle 3, 0 \rangle$	<b>a b a</b> b a b $\dots$
$\langle 1, 2 \rangle$	<b>b a a</b> b a a $\dots$
$\langle 2, 2 \rangle$	<b>b a a</b> b a a $\dots$
$\langle 3, 1 \rangle$	<b>b a b</b> a b a $\dots$



# Longest common prefix and distinguishing prefix arrays

$$W = \{v_1, v_2, v_3\} = \{aab, aab, ab\}$$

$SA_W$	$\text{suf}(W)$	$LCP_W$	$DP_W$
$\langle 1, 0 \rangle$	<b>a a b</b> a a b $\dots$	—	—
$\langle 2, 0 \rangle$	<b>a a b</b> a a b $\dots$	$\infty$	$\infty$
$\langle 1, 1 \rangle$	<b>a b a</b> a b a $\dots$	1	2
$\langle 2, 1 \rangle$	<b>a b a</b> a b a $\dots$	$\infty$	$\infty$
$\langle 3, 0 \rangle$	<b>a b a</b> b a b $\dots$	3	4
$\langle 1, 2 \rangle$	<b>b a a</b> b a a $\dots$	0	1
$\langle 2, 2 \rangle$	<b>b a a</b> b a a $\dots$	$\infty$	$\infty$
$\langle 3, 1 \rangle$	<b>b a b</b> a b a $\dots$	2	3

# Burrows-Wheeler transform

$$W = \{v_1, v_2, v_3\} = \{aab, aab, ab\}$$

$SA_W$	$\text{suf}(W)$	$LCP_W$	$DP_W$	$BWT(W)$
$\langle 1, 0 \rangle$	<b>a a b</b> a a b $\dots$	—	—	b
$\langle 2, 0 \rangle$	<b>a a b</b> a a b $\dots$	$\infty$	$\infty$	b
$\langle 1, 1 \rangle$	<b>a b a</b> a b a $\dots$	1	2	a
$\langle 2, 1 \rangle$	<b>a b a</b> a b a $\dots$	$\infty$	$\infty$	a
$\langle 3, 0 \rangle$	<b>a b a</b> b a b $\dots$	3	4	b
$\langle 1, 2 \rangle$	<b>b a a</b> b a a $\dots$	0	1	a
$\langle 2, 2 \rangle$	<b>b a a</b> b a a $\dots$	$\infty$	$\infty$	a
$\langle 3, 1 \rangle$	<b>b a b</b> a b a $\dots$	2	3	a

# Cyclic equivalence

## Cyclic equivalence

Two multisets of words  $V$  and  $W$  are cyclically equivalent if

$$\text{suf}(V) = \text{suf}(W).$$

## Example multiset

$$W = \{\{aab, aab, ab\}\}$$

*aabaabaabaab...*

*abaabaabaaba...*

*baabaabaabaa...*

*aabaabaabaab...*

*abaabaabaaba...*

*baabaabaabaa...*

*abababababab...*

*babababababa...*

## Equivalence class

$\{\{aab, aab, ab\}\}$

$\{\{aba, aab, ab\}\}$

$\{\{baa, aab, ab\}\}$

$\{\{aab, aba, ab\}\}$

$\vdots$

$\{\{aabaab, ab\}\}$

$\{\{abaaba, ab\}\}$

$\{\{baabaa, ab\}\}$

$\{\{aab, aab, ba\}\}$

$\{\{aba, aab, ba\}\}$

$\{\{baa, aab, ba\}\}$

$\{\{aab, aba, ba\}\}$

$\vdots$

$\{\{aabaab, ba\}\}$

$\{\{abaaba, ba\}\}$

$\{\{baabaa, ba\}\}$



## Lemma

Let  $W = \{\{w_i\}_{i=1}^s$  be a multiset of cyclic words. Then:

- 1 There exists a set of cyclic words  $V = \{v_i\}_{i=1}^t$  such that  $\text{suf}(W) = \text{suf}(V)$ .
- 2 There exists a multiset of primitive cyclic words  $U = \{\{u_i\}_{i=1}^p$  such that  $\text{suf}(W) = \text{suf}(U)$ .

## Lemma

Let  $W = \{\{w_i\}_{i=1}^s$  be a multiset of cyclic words. Then:

- 1 There exists a set of cyclic words  $V = \{v_i\}_{i=1}^t$  such that  $\text{suf}(W) = \text{suf}(V)$ .
- 2 There exists a multiset of primitive cyclic words  $U = \{\{u_i\}_{i=1}^p$  such that  $\text{suf}(W) = \text{suf}(U)$ .

## Remark

If two multisets of words  $V$  and  $W$  are cyclically equivalent, then  $\text{LCP}_V = \text{LCP}_W$ ,  $\text{DP}_V = \text{DP}_W$  and  $\text{BWT}_V = \text{BWT}_W$ .

## Lemma

Let  $W = \{\{w_i\}_{i=1}^s\}$  be a multiset of cyclic words. Then:

- 1 There exists a set of cyclic words  $V = \{v_i\}_{i=1}^t$  such that  $\text{suf}(W) = \text{suf}(V)$ .
- 2 There exists a multiset of primitive cyclic words  $U = \{\{u_i\}_{i=1}^p\}$  such that  $\text{suf}(W) = \text{suf}(U)$ .

## Remark

If two multisets of words  $V$  and  $W$  are cyclically equivalent, then  $\text{LCP}_V = \text{LCP}_W$ ,  $\text{DP}_V = \text{DP}_W$  and  $\text{BWT}_V = \text{BWT}_W$ .

## Theorem (Mantaci, Restivo, Rosone, Sciortino – 2007)

The mapping from a word  $v$  to the cyclical equivalence class of  $\text{IBWT}(v)$  is a bijection

## Irreducible values

A value  $LCP_W[i]$  (respectively  $DP_W[i]$ ) is irreducible if  $BWT_W[i] \neq BWT_W[i - 1]$ .

$SA_W$	$\text{suf}(W)$	$LCP_W$	$DP_W$	$BWT(W)$
$\langle 1, 0 \rangle$	<b>a a b</b> a a b ...	—	—	b
$\langle 2, 0 \rangle$	<b>a a b</b> a a b ...	$\infty$	$\infty$	b
$\langle 1, 1 \rangle$	<b>a b a</b> a b a ...	<b>1</b>	<b>2</b>	a
$\langle 2, 1 \rangle$	<b>a b a</b> a b a ...	$\infty$	$\infty$	a
$\langle 3, 0 \rangle$	<b>a b a</b> b a b ...	<b>3</b>	<b>4</b>	b
$\langle 1, 2 \rangle$	<b>b a a</b> b a a ...	<b>0</b>	<b>1</b>	a
$\langle 2, 2 \rangle$	<b>b a a</b> b a a ...	$\infty$	$\infty$	a
$\langle 3, 1 \rangle$	<b>b a b</b> a b a ...	2	3	a

## Irreducible values

- $W = \{w_i\}_{i=1}^s$  and  $\sum_{i=1}^s |w_i| = n$
- $\Sigma_{\text{lcp}}(W)$  – sum of irreducible LCP values
- $\Sigma_{\text{idp}}(W)$  – sum of distinguishing prefixes lengths

## Theorem (Kärkkäinen, Manzini, Puglisi – 2009)

$$\Sigma_{\text{lcp}}(W) = O(n \log n)$$

# New upper bounds for the sum of irreducible lcp values

## Theorem

For any multiset  $W$  of words of total length  $n > 0$ , we have

$$\Sigma_{\text{ilcp}}(W) \leq \Sigma_{\text{idp}}(W) \leq n \lceil \lg n \rceil - 2^{\lceil \lg n \rceil} + 1$$

## Theorem

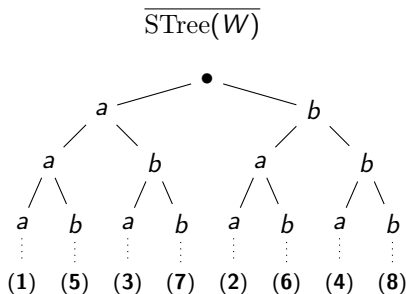
For any multiset  $W$  of words of total length  $n > 0$  such that  $\text{BWT}(W)$  has  $r$  runs, we have

$$\Sigma_{\text{ilcp}}(W) + r - 1 = \Sigma_{\text{idp}}(W) \leq n \lceil \lg r \rceil - 2^{\lceil \lg r \rceil} + 1$$

# Reverse suffixes

$$W = \{a, aab, abb, b\}$$

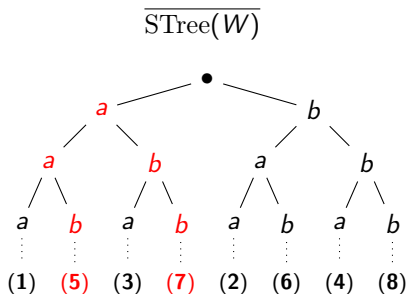
- (1) *aaaaaa...*
- (2) *baabaa...*
- (3) *abaaba...*
- (4) *bbabba...*
- (5) *aabaab...*
- (6) *babbab...*
- (7) *abbabb...*
- (8) *bbbbbb...*



# Reverse suffixes

$$W = \{a, aab, abb, b\}$$

- (1) *aaaaaa...*
- (2) *baabaa...*
- (3) *abaaba...*
- (4) *bbabba...*
- (5) *aabaab...*
- (6) *babbab...*
- (7) *abbabb...*
- (8) *bbbbbb...*



Dispersion pair with respect to  $\leq$

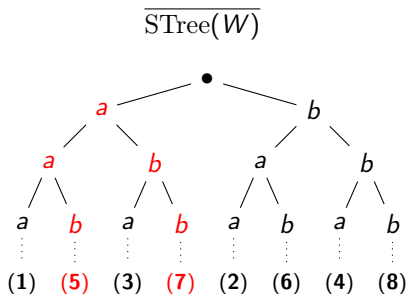
The pair of leaves  $(u, v)$  of an ordered tree is called a dispersion pair if  $u < v$  the subtree rooted at their nearest common ancestor contains no leaf  $w$  such that  $u < w < v$ .



# Reverse suffixes

$$W = \{a, aab, abb, b\}$$

- (1) *aaaaaa...*
- (2) *baabaa...*
- (3) *abaaba...*
- (4) *bbabba...*
- (5) *aabaab...*
- (6) *babbab...*
- (7) *abbabb...*
- (8) *bbbbbb...*



## Lemma

$$D(\overline{\text{STree}(W)}, \leq_W) = \Sigma \text{idp}(W)$$

$D(T, \leq)$  – the number of dispersal pairs in  $T$

$$d(1) = 0$$

$$d(n) = \max_{i \in [1.. \lfloor n/2 \rfloor]} d(n, i) \quad \text{when } n > 1$$

$$d(n, k) = d(k) + d(n - k) + \min\{2k, n - 1\}$$

where  $n > 0$  and  $k \in [0.. \lfloor n/2 \rfloor]$ .

## Lemma

$d(n) = \max\{D(T, \leq)\}$ , where the maximum is taken over any rooted tree  $T$  with  $n$  leaves and any total order  $\leq$  on its leaves.

## Lemma

$$d(n) = n \lceil \lg n \rceil - 2^{\lceil \lg n \rceil} + 1$$

## Theorem

For any multiset  $W$  of words of total length  $n > 0$ , we have

$$\sum \text{ilcp}(W) \leq \sum \text{idp}(W) \leq d(n) \leq n \lceil \lg n \rceil - 2^{\lceil \lg n \rceil} + 1$$

## Lemma

If  $\text{BWT}(W)$  has  $r$  runs, then  $\left| D_u(\overline{\text{STree}}(W), \leq_W) \right| < r$  for every vertex  $u$  in  $\overline{\text{STree}}(W)$ .

$$d_r(1) = 0$$

$$d_r(n) = \max_{i \in [1.. \lfloor n/2 \rfloor]} d_r(n, i) \quad \text{when } n > 1$$

$$d_r(n, k) = d_r(k) + d_r(n - k) + \min\{2k, n - 1, r - 1\}$$

where  $r > 0$ ,  $n > 0$  and  $k \in [0.. \lfloor n/2 \rfloor]$ .

## Lemma

$d_r(n) = \max\{d(T, \leq)\}$ , where the maximum is taken over any rooted tree  $T$  with  $n$  leaves and any total order  $\leq$  on its leaves s.t.  $|D_u(T, \leq)| < r$  for every vertex  $u \in T$ .

## Lemma

For any  $2 \leq r \leq n$  we have  $d_r(n) \leq n \lceil \lg r \rceil - 2^{\lceil \lg r \rceil} + 1$ .

## Theorem

For any multiset  $W$  of words of total length  $n > 0$  such that  $\text{BWT}(W)$  has  $r$  runs, we have

$$\sum \text{ilcp}(W) + r - 1 = \sum \text{idp}(W) \leq d_r(n) \leq n \lceil \lg r \rceil - 2^{\lceil \lg r \rceil} + 1$$

## De Bruijn set

A set of words  $W = \{w_i\}_{i=1}^t$  over the alphabet  $\mathcal{A}$  is a de Bruijn set of order  $k$  if  $\sum_{i=1}^t |w_i| = 2^k$  and every word  $v \in \mathcal{A}^k$  is a prefix of exactly one word in  $\text{suf}(W)$ .

$$W = \{a, aab, abb, b\}$$

$$\sum_{i=1}^4 |w_i| = 2^3$$

$$W_{\langle 1,1 \rangle} = \mathbf{a} a a a a \dots$$

$$W_{\langle 2,1 \rangle} = \mathbf{a a b} a a b \dots$$

$$W_{\langle 2,3 \rangle} = \mathbf{a b a} a b a \dots$$

$$W_{\langle 3,1 \rangle} = \mathbf{a b b} a b b \dots$$

$$W_{\langle 2,3 \rangle} = \mathbf{b a a} b a a \dots$$

$$W_{\langle 3,2 \rangle} = \mathbf{b a b} b a b \dots$$

$$W_{\langle 4,3 \rangle} = \mathbf{b b a} b b a \dots$$

$$W_{\langle 4,1 \rangle} = \mathbf{b b b b b} \dots$$

## De Bruijn set

A set of words  $W = \{w_i\}_{i=1}^t$  over the alphabet  $\mathcal{A}$  is a de Bruijn set of order  $k$  if  $\sum_{i=1}^t |w_i| = 2^k$  and every word  $v \in \mathcal{A}^k$  is a prefix of exactly one word in  $\text{suf}(W)$ .

$$W = \{a, aab, abb, b\}$$

$$\sum_{i=1}^4 |w_i| = 2^3$$

$$W_{\langle 1,1 \rangle} = \mathbf{a} a a a a \dots$$

$$W_{\langle 2,3 \rangle} = \mathbf{b} a a b a a \dots$$

$$W_{\langle 2,1 \rangle} = \mathbf{a a b} a a b \dots$$

$$W_{\langle 3,2 \rangle} = \mathbf{b a b} b a b \dots$$

$$W_{\langle 2,3 \rangle} = \mathbf{a b a} a b a \dots$$

$$W_{\langle 4,3 \rangle} = \mathbf{b b a} b b a \dots$$

$$W_{\langle 3,1 \rangle} = \mathbf{a b b} a b b \dots$$

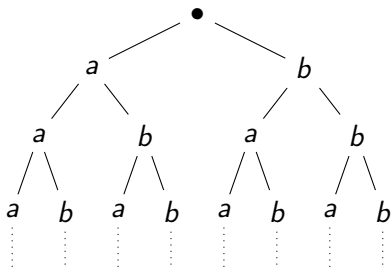
$$W_{\langle 4,1 \rangle} = \mathbf{b b b b b} \dots$$

## Lemma (Higgins – 2012)

For  $k \geq 1$ , and any  $u \in U_k = \{ab, ba\}^{2^{k-1}}$ ,  $W = \text{IBWT}(u)$  is a de Bruijn set of order  $k$ .



$$W_k = \text{IBWT}((ab)^{k-1})$$



$$W_3 = \{a, aab, abb, b\}$$

### Theorem

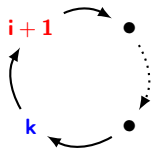
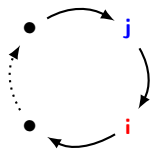
$$\Sigma \text{ilcp}(W_k) = \sum_{i=1}^{k-1} i2^i = k \cdot 2^k - 2^{k+1} + 2$$

## Theorem

For any  $k \geq 1$ , there exists a word  $w$  of length  $n = 2^k$  such that  $\Sigma \text{idp}(w) = n \log n - O(n)$ .

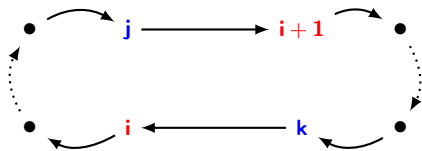
## Theorem

For any  $k \geq 1$ , there exists a word  $w$  of length  $n = 2^k$  such that  $\Sigma \text{idp}(w) = n \log n - O(n)$ .

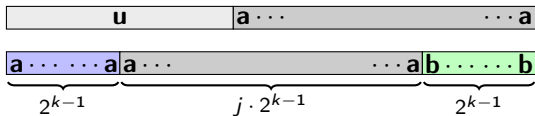
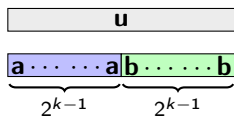


## Theorem

For any  $k \geq 1$ , there exists a word  $w$  of length  $n = 2^k$  such that  $\Sigma \text{idp}(w) = n \log n - O(n)$ .

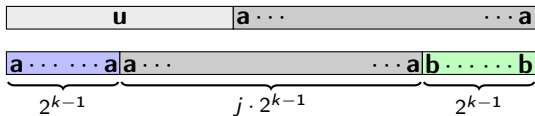
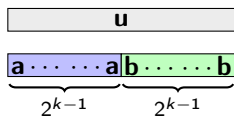


$$W_{k,j} = \text{IBWT} \left( \underbrace{(ab)^{2^{k-1}}}_{u} a^{j2^{k-1}} \right)$$



# $n \log r$ lower bound for $\Sigma \text{idp}$

$$W_{k,j} = \text{IBWT} \left( \underbrace{(ab)^{2^{k-1}}}_{u} a^{j2^{k-1}} \right)$$



$$W_{k,j} = \bigcup_{i=1}^{j+1} S_{i,k,j}$$

$$S_{i,k,j} = \begin{cases} a^i b a^j \{a, b a^j\}^{k-1} & : i \leq j \\ a^{j+1} \{a, b a^j\}^{k-1} & : i > j \end{cases}$$

## Lemma

For  $j \geq 1$  we have  $\Sigma \text{lcp}(W_{k,j}) = (j+2)k2^{k-1} - 2^{k+1} + j + 1$ .

## Theorem

For any  $r = 2^k + 1$ ,  $k \geq 1$ , and  $n \geq r$  such that  $2^{k-1} | n$ , there exists a word  $w$  of length  $n$  such that  $\text{BWT}(w)$  contains  $r - o(r)$  runs and  $\Sigma \text{idp}(w) = n \log r - O(n)$ .

- New upper bound for  $\Sigma_{ilcp}$  and  $\Sigma_{idp}$
- New upper bound for  $\Sigma_{ilcp}$  and  $\Sigma_{idp}$  related to the number of runs in BWT
- Lower bounds for  $\Sigma_{ilcp}$  and  $\Sigma_{idp}$  matching upper bounds
- Lower bounds for  $\Sigma_{ilcp}$  and  $\Sigma_{idp}$  – single word case





**THANK YOU**



**KIITOS**



**DZIĘKUJĘ**



**GRAZIE**