

- LCA-esikäsitelyä voidaan käyttää esim. palindromien etsimiseen.
  - Palindromi on jono, joka on oma käänteisjono. Esimerkiksi SAIPPUAKAUPPIAS on palindromi.
  - Kuvataan, miten lineaarisessa ajassa löydetään pisin parittoman mittainen palindromi, joka esiintyy jonossa  $S$ . Pisin parillisen mittainen palindromi löydetään vastaavasti.
  - Olkoon  $S^R$  jonon  $S$  käänteisjono. Jos  $S[i - k \dots i + k]$  on pisin palindromi, jonka keskikohta on  $i$ , niin loppuosien  $S[i \dots n]$  ja  $S^R[n - i - 1 \dots n]$  pisin yhteinen alkuosa on  $S[i \dots i + k] = S^R[n - i - 1 \dots n - i - 1 + k]$ .
  - Se löydetään vakioajassa jonojen  $S$  ja  $S^R$  yleistetystä loppuosapuusta, joka on LCA-esikäsitely. Vastaava kysely voidaan tehdä kaikille mahdollisille keskikohdille ajassa  $\mathcal{O}(n)$ .

## 4.4 Loppuosataulukko

U. Manber & G. Myers: Suffix Arrays: A New Method for On-line String Searches. *SIAM Journal on Computing* 22(5):935–948, 1993.

- Merkkijonon  $S$  loppuosataulukko on loppuosien joukon  $S_{[0,n]}$  järjestetty taulukko.
- Täsmällisemmin,  $S$ :n loppuosataulukko on taulukko  $SA[0 \dots n]$ , joka sisältää joukon  $[0, n]$  permutaation ja toteuttaa ehdon  $S_{SA[0]} < S_{SA[1]} < \dots < S_{SA[n]}$ .
- Hahmonsovitusta eli alkuosahaku hahmolla  $P$  voidaan tehdä ajassa  $\mathcal{O}(|P| + \log n + occ)$  binäärihaulla olettaen, että myös tehostetussa binäärihaussa tarvittavat *LLCP* ja *RLCP* taulukot on annettu.
- Monissa loppuosataulukon sovelluksissa on keskeisessä asemassa toinen LCP-taulukko:  $LCP[i] = lcp(S_{SA[i]}, S_{SA[i+1]})$ .
- Samoin usein esiintyvä taulukko on loppuosataulukon käänteistaulukko:  $SA^{-1}[i] = j$  jos ja vain jos  $SA[j] = i$ .
- $SA^{-1}[i]$  voidaan myös tulkita loppuosan  $S_i$  järjestysnumeroksi.

**Esimerkki 4.4.1** Jonon  $S = \text{banana}$  loppuosataulukko, LCP-taulukko ja loppuosataulukon käänteistaulukko.

	$SA$	$LCP$		$SA^{-1}$	
0	6	0		4	banana
1	5	1	a	3	anana
2	3	3	ana	6	nana
3	1	0	anana	2	ana
4	0	0	banana	5	na
5	4	2	na	1	a
6	2	0	nana	0	

### LCP-taulukon laskeminen

T. Kasai, G. Lee, H. Arimura, S. Arikawa & K. Park: Linear-Time Longest-Common-Prefix Computation in Suffix Arrays and Its Applications. Teoksessa *Proc. CPM 2001*, LNCS 2089, sivut 181–192, 2001.

- Käänteistaulukko  $SA^{-1}$  on helppo laskea lineaarisessa ajassa loppuosataulukosta  $SA$ .
- Sen avulla myös LCP-taulukko voidaan helposti laskea lineaarisessa ajassa

**Lemma 4.4.2** *Kaikilla  $0 \leq i < n$ ,  $LCP[SA^{-1}[i+1]] \geq LCP[SA^{-1}[i]] - 1$ .*

*Todistus.*

- Oletetaan, että  $S_i$  ja  $S_j$  ovat aakkosjärjestyksessä peräkkäin eli  $LCP[SA^{-1}[i]] = lcp(S_i, S_j) = \ell$ .
- Jos  $\ell > 0$ , niin jollakin merkillä  $a$ ,  $S_i = aS_{i+1}$  ja  $S_j = aS_{j+1}$  ja siten  $S_{i+1} < S_{j+1}$  ja  $lcp(S_{i+1}, S_{j+1}) = \ell - 1$ .
- Tällöin  $LCP[SA^{-1}[i+1]] \geq lcp(S_{i+1}, S_{j+1}) = \ell - 1$ .

□

Seuraava algoritmi laskee LCP-taulukon järjestyksessä  $LCP[SA^{-1}[0]]$ ,  $LCP[SA^{-1}[1]]$ ,  $LCP[SA^{-1}[2]]$ , ...

**Algoritmi 4.4.3** *LCP-taulukon muodostaminen*

Syöte: Merkkijono  $S = S[0 \dots n]$ , loppuosataulukko  $SA[0 \dots n]$   
ja sen käänteistaulukko  $SA^{-1}[0 \dots n]$ .

Tuloste: LCP-taulukko  $LCP[0 \dots n]$ .

```

(1)  $\ell := 0$ ;
(2) for  $i := 0$  to  $n - 1$  do
(3)    $k := SA^{-1}[i]$ ;           (*  $i = SA[k]$  *)
(4)   if  $k = n$  then  $LCP[k] := 0$ ;
(5)   else
(6)      $j := SA[k + 1]$ ;
(7)      $\ell := \ell + lcp(S_{i+\ell}, S_{j+\ell})$ ;
(8)      $LCP[k] := \ell$ ;           (*  $LCP[k] := lcp(S_i, S_j)$  *)
(9)   if  $\ell > 0$  then  $\ell := \ell - 1$ ;

```

**Lause 4.4.4** *Annettuna merkkijono  $S[0 \dots n]$  ja sen loppuosataulukko  $SA[0 \dots n]$ , vastaava LCP-taulukko  $LCP[0 \dots n]$  voidaan laskea ajassa  $\mathcal{O}(n)$ .*

*Todistus.*

- Algoritmi toimii selvästi lineaarisessa ajassa lukuunottamatta rivin 7 lcp-arvon laskemista.
- Jokainen ylimääräinen merkkivertailu rivillä 7 kasvattaa  $\ell$ :ää.
- Koska  $\ell$  pienenee enintään yhdellä per kierros, eikä koskaan kasva suuremmaksi kuin  $n$ , rivillä 7 tehdään yhteensä  $\mathcal{O}(n)$  merkkivertailua.

□

## 4.5 Loppuosataulukon muodostaminen

S. J. Puglisi, W. F. Smyth & A. Turpin. A taxonomy of suffix array construction algorithms, *ACM Computing Surveys*, 2007.

- Tarkastellaan nyt algoritmeja loppuosataulukon muodostamiseen.
- Tehtävänä on siis järjestää merkkijonon  $S[0, n]$  loppuosat aakkosjärjestykseen.

- Normaali merkkijonojen järjestämialgoritmi vie pahimmassa tapauksessa ajan  $\Omega(n^2)$ .
- Loppuosataulukko voidaan muodostaa laskemalla ensin loppuosapuun ja siitä loppuosataulukko. Tämä toimii myös toiseen suuntaan: loppuosapuun on helppo muodostaa lineaarisessa ajassa loppuosataulukosta ja LCP-taulukosta.
- Esitetään nyt tehokkaita algoritmeja, joilla loppuosataulukko voidaan muodostaa suoraan ilman loppuosapuuta.

### 4.5.1 Karp-Miller-Rosenberg nimentätekniikka

R.M. Karp & R.E. Miller & A.L. Rosenberg: Rapid Identification of Repeated Patterns in Strings Trees and Arrays. *Proceedings of the 4th Symposium on Theory of Computing*, sivut 125–136, 1972.

- Merkitään  $S_i^\ell = S[i \dots \min\{i + \ell, n\}]$ .
- Muodostetaan taulukot  $Name_\ell[0, n]$ , jotka sisältävät kokonaislukuja väliltä  $[0, n]$  ja toteuttavat kaikilla  $i, j \in [0, n]$

$$Name_\ell[i] \leq Name_\ell[j] \iff S_i^\ell \leq S_j^\ell$$

- Taulukko  $Name_\ell[0, n]$  antaa siis  $\ell$ :n pituisten osajonojen ja  $\ell$ :ää lyhyempien loppuosien nimennän, joka säilyttää jonojen keskinäisen järjestyksen.
- Erityisesti, kun  $\ell \geq n$ , saadaan kaikkien loppuosien nimentä. Siis  $Name_\ell[i] = SA^{-1}[i]$ , mistä saadaan helposti loppuosataulukko.
- Yllä oleva nimennän määritelmä ei ole yksikäsitteinen, mutta nimeksi voi valita

$$Name_\ell[i] = |\{j \in [0, n] \mid S_j^\ell < S_i^\ell\}|$$

- Taulukko  $Name_1[0, n]$  voidaan laskea järjestämällä yhden pituiset jonot, siis jonon  $S$  merkit, ajassa  $\mathcal{O}(n \log n)$ .
- Taulukko  $Name_{2\ell}[0, n]$  voidaan muodostaa taulukosta  $Name_\ell[0, n]$  järjestämällä parit  $(Name_\ell[i], Name_\ell[i + \ell])$ . Tämä voidaan tehdä  $\mathcal{O}(n)$  ajassa käyttämällä lopusta alkuun kantalukujärjestämistä.

- Loppuosataulukko (tai sen käänteistaulukko) saadaan toistamalla kahdentamisaskelta  $\lceil \log n \rceil$  kertaa eli laskemalla  $Name_1, Name_2, Name_4, Name_8, \dots$

**Esimerkki 4.5.1** Käänteisen loppuosataulukon laskenta jonolle *banana*.

$Name_1:$	4	b	$Name_2:$	4	ba	$Name_4:$	4	bana
	1	a		2	an		3	anan
	5	n		5	na		6	nana
	1	a		2	an		2	ana
	5	n		5	na		5	na
	1	a		1	a		1	a
	0			0			0	

$Name_8 = SA^{-1}:$	4	banana
	3	anana
	6	nana
	2	ana
	5	na
	1	a
	0	

**Lause 4.5.2** Merkkijonon  $S[0, n]$  loppuosataulukko voidaan laskea ajassa  $\mathcal{O}(n \log n)$ .

- Karp-Miller-Rosenberg nimentäteknikkaa voidaan käyttää myös muihin ongelmiin.
- Loppuosataulukon muodostamiseen sitä sovelsivat ensin Manber ja Myers.
- Käytännössä parhaan version tästä algoritmista ovat esittäneet Larsson ja Sadakane, joiden algoritmi sisältää joitakin merkkijonopikajärjestämisen piirteitä.

### 4.5.2 Osarekursio

J. Kärkkäinen & P. Sanders & S. Burkhardt: Linear Work Suffix Array Construction. *Journal of the ACM* 53(6):918–936, 2006.

- Esitetään nyt algoritmi, joka laskee loppuosataulukon ajassa  $\mathcal{O}(n)$ , kun aakkosto on  $[1, n]$ .
- Jos aakkosto on joku muu kuin  $[1, n]$ , voidaan kukin merkki  $S[i]$  korvata nimellä  $Name_1[i]$ , jolloin päästään aakkostoon  $[1, n]$  muuttamatta loppuosien keskinäistä järjestystä.
- Siten millä tahansa aakkostolla jonon  $S$  loppuosat voidaan järjestää samassa ajassa kuin  $S$ :n merkit.
- Perusrakenne on seuraava osarekursio:
  0. Jaa loppuosat kahteen ryhmään  $S_A$  ja  $S_B$ .
  1. Järjestä  $S_A$ . Tämä tapahtuu palauttamalla ongelma loppuosataulukon muodotukseen jonolle, jonka pituus on  $|A|$ .
  2. Järjestä  $S_B$  käyttäen  $S_A$ :n järjestystä.
  3. Lomita  $S_A$  ja  $S_B$ .
- Rekursio poislukien kaikki pystytään tekemään lineaarisessa ajassa.
- Lisäksi pätee  $|A| \leq 2n/3$ . Siten aikavaatimus on  $T(n) = \mathcal{O}(n) + T(2n/3) = \mathcal{O}(n)$ .
- Olennaista on, että  $S_B$ :tä ei järjestetä rekursiivisesti.
- Samankaltainen osarekursiivinen rakenne esiintyy kaikissa loppuosataulukon ja -puun muodostusalgoritmeissa, jotka toimivat lineaarisessa ajassa aakkostolla  $[1, n]$ .
- Syötteenä on siis jono  $S[0, n)$  aakkostossa  $[1, n]$ . Lisäksi oletetaan, että  $S[n, n+3) = 000$ .

#### Esimerkki 4.5.3

$$S[0, n) = \begin{array}{cccccccccccc} & 0 & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 \\ \text{y} & \text{a} & \text{b} & \text{b} & \text{a} & \text{d} & \text{a} & \text{b} & \text{b} & \text{a} & \text{d} & \text{o} & \end{array}$$

**Askel 0: Muodosta  $S_A$ .**

- Määritellään kaikilla  $k = 0, 1, 2$ ,  $C_k = \{i \in [0, n] \mid i \bmod 3 = k\}$ .  
Olkoon  $A = C_1 \cup C_2$  ja  $B = C_0$ .

**Esimerkki 4.5.4**  $B = C_0 = \{0, 3, 6, 9, 12\}$ ,  $C_1 = \{1, 4, 7, 10\}$ ,  $C_2 = \{2, 5, 8, 11\}$ , eli  $A = \{1, 4, 7, 10, 2, 5, 8, 11\}$ .

**Askel 1: Järjestä  $S_A$ .**

- Muodosta jonot

$$R_k = \left[ S[k \dots k+3] \right] \left[ S[k+3 \dots k+6] \right] \dots \left[ S[\max C_k \dots \max C_k + 3] \right]$$

$k = 1, 2$ , joiden merkit ovat kolmen merkin pituisia  $S$ :n osajonoja.

- Olkoon  $R$  jonojen  $R_1$  ja  $R_2$  katenaatio. Sen loppuosat vastaavat  $S_A$ :n loppuosia.

**Esimerkki 4.5.5**  $R = [abb][ada][bba][doo][bba][dab][bad][ooo]$ .

- Korvataan  $R$ :n merkit KMR-tekniikan mukaisilla nimillä. Tämä vaatii  $R$ :n merkkien järjestämistä, mikä voidaan tehdä lopusta alkuun kantalukujärjestämisellä lineaarisessa ajassa. Olkoon  $R'$  muunnettu jono.
- Muodostetaan rekursiivisella kutsulla  $R'$ :n loppuosataulukko, mistä saadaan  $S_A$ :n järjestys. Rekursiota ei tarvita, jos kaikki  $R'$ :n merkit ovat erilaiset.

**Esimerkki 4.5.6**  $R' = (1, 2, 4, 7, 4, 6, 3, 8)$ ,  
 $SA_{R'}^{-1} = (1, 2, 5, 7, 4, 6, 3, 8, 0)$  ja  $SA_{R'} = (8, 0, 1, 6, 4, 2, 5, 3, 7)$ .

- Olkoon  $rank(i)$ ,  $i \in A$ , loppuosan  $S_i$  järjestysluku joukossa  $S_A$ . Lisäksi määritellään  $rank(n+1) = rank(n+2) = 0$ . Kun  $i \in C_0$ ,  $rank(i) = \perp$ .

**Esimerkki 4.5.7**

$rank(S_i)$	$i$	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14
		$\perp$	1	4	$\perp$	2	6	$\perp$	5	3	$\perp$	7	8	$\perp$	0	0

**Askel 2: Järjestä  $S_B$ .**

- Muodostetaan jokaiselle  $S_i \in S_B$  pari  $(t_i, \text{rank}(S_{i+1}))$ . Huomaa, että  $\text{rank}(S_{i+1})$  on aina määritelty.
- Selvästi kaikilla  $i, j \in B$ ,

$$S_i \leq S_j \iff (S[i], \text{rank}(S_{i+1})) \leq (S[j], \text{rank}(S_{j+1})) .$$

- Parit  $(S[i], \text{rank}(S_{i+1}))$  järjestetään lopusta alkuun kantalukujärjestämisellä.

**Esimerkki 4.5.8**  $S_{12} < S_6 < S_9 < S_3 < S_0$  koska  $(0, 0) < (a, 5) < (a, 7) < (b, 2) < (y, 1)$ .

**Askel 3: Lomita.**

- Lomitetaan  $S_A$  ja  $S_B$  normaalilla vertailuihin perustuvalla lomitusalgoritmeilla.
- Kahden loppuosan  $S_i \in S_A$  ja  $S_j \in S_B$  vertailussa erotetaan kaksi tapausta:

$$i \in C_1 : \quad S_i \leq S_j \iff (S[i], \text{rank}(i+1)) \leq (S[j], \text{rank}(j+1))$$

$$i \in C_2 : \quad S_i \leq S_j \iff (S[i], S[i+1], \text{rank}(i+2)) \leq (S[j], S[j+1], \text{rank}(j+2))$$

**Esimerkki 4.5.9**  $S_1 < S_6$  koska  $(a, 4) < (a, 5)$  ja  $S_3 < S_8$  koska  $(b, a, 6) < (b, a, 7)$ .

**Toteutus**

- Alla on algoritmin täydellinen C++-kielinen toteutus.
- Toteutus ei ole tehokkain mahdollinen; pyrkimys on lyhyteen (50 riviä poislukien tyhjät rivit, kommentit, ja rivit, joilla on vain aaltosulku).

```
inline bool leq(int a1, int a2, int b1, int b2) // lexicographic order
{ return(a1 < b1 || a1 == b1 && a2 <= b2); } // for pairs
inline bool leq(int a1, int a2, int a3, int b1, int b2, int b3)
{ return(a1 < b1 || a1 == b1 && leq(a2,a3, b2,b3)); } // and triples
```



```

// stably sort a[0..n-1] to b[0..n-1] with keys in 0..K from r
static void radixPass(int* a, int* b, int* r, int n, int K)
{ // count occurrences
  int* c = new int[K + 1]; // counter array
  for (int i = 0; i <= K; i++) c[i] = 0; // reset counters
  for (int i = 0; i < n; i++) c[r[a[i]]]++; // count occurrences
  for (int i = 0, sum = 0; i <= K; i++) // exclusive prefix sums
  { int t = c[i]; c[i] = sum; sum += t; }
  for (int i = 0; i < n; i++) b[c[r[a[i]]]++] = a[i]; // sort
  delete [] c;
}

// find the suffix array SA of T[0..n-1] in {1..K}^n
// require T[n]=T[n+1]=T[n+2]=0, n>=2
void suffixArray(int* T, int* SA, int n, int K) {
  int n0=(n+2)/3, n1=(n+1)/3, n2=n/3, n02=n0+n2;
  int* R = new int[n02 + 3]; R[n02]= R[n02+1]= R[n02+2]=0;
  int* SA12 = new int[n02 + 3]; SA12[n02]=SA12[n02+1]=SA12[n02+2]=0;
  int* R0 = new int[n0];
  int* SA0 = new int[n0];

  //***** Step 0: Construct sample *****
  // generate positions of mod 1 and mod 2 suffixes
  // the "+(n0-n1)" adds a dummy mod 1 suffix if n%3 == 1
  for (int i=0, j=0; i < n+(n0-n1); i++) if (i%3 != 0) R[j++] = i;

  //***** Step 1: Sort sample suffixes *****
  // lsb radix sort the mod 1 and mod 2 triples
  radixPass(R , SA12, T+2, n02, K);
  radixPass(SA12, R , T+1, n02, K);
  radixPass(R , SA12, T , n02, K);

  // find lexicographic names of triples and
  // write them to correct places in R
  int name = 0, c0 = -1, c1 = -1, c2 = -1;
  for (int i = 0; i < n02; i++) {
    if (T[SA12[i]] != c0 || T[SA12[i]+1] != c1 || T[SA12[i]+2] != c2)
    { name++; c0 = T[SA12[i]]; c1 = T[SA12[i]+1]; c2 = T[SA12[i]+2]; }
    if (SA12[i] % 3 == 1) { R[SA12[i]/3] = name; } // write to R1
    else { R[SA12[i]/3 + n0] = name; } // write to R2
  }

  // recurse if names are not yet unique
  if (name < n02) {
    suffixArray(R, SA12, n02, name);
    // store unique names in R using the suffix array
    for (int i = 0; i < n02; i++) R[SA12[i]] = i + 1;
  } else // generate the suffix array of R directly
    for (int i = 0; i < n02; i++) SA12[R[i] - 1] = i;
}

```

```

//***** Step 2: Sort nonsample suffixes *****
// stably sort the mod 0 suffixes from SA12 by their first character
for (int i=0, j=0; i < n02; i++) if (SA12[i] < n0) R0[j++] = 3*SA12[i];
radixPass(R0, SA0, T, n0, K);

//***** Step 3: Merge *****
// merge sorted SA0 suffixes and sorted SA12 suffixes
for (int p=0, t=n0-n1, k=0; k < n; k++) {
#define GetI() (SA12[t] < n0 ? SA12[t] * 3 + 1 : (SA12[t] - n0) * 3 + 2)
  int i = GetI(); // pos of current offset 12 suffix
  int j = SA0[p]; // pos of current offset 0 suffix
  if (SA12[t] < n0 ? // different compares for mod 1 and mod 2 suffixes
      leq(T[i], R[SA12[t] + n0], T[j], R[j/3]) :
      leq(T[i], T[i+1], R[SA12[t]-n0+1], T[j], T[j+1], R[j/3+n0]))
  {
    // suffix from SA12 is smaller
    SA[k] = i; t++;
    if (t == n02) // done --- only SA0 suffixes left
      for (k++; p < n0; p++, k++) SA[k] = SA0[p];
  } else { // suffix from SA0 is smaller
    SA[k] = j; p++;
    if (p == n0) // done --- only SA12 suffixes left
      for (k++; t < n02; t++, k++) SA[k] = GetI();
  }
}
delete [] R; delete [] SA12; delete [] SA0; delete [] R0;
}

```

## 4.6 Loppuosataulukon sovelluksia

Useimmat loppuosapuun sovellukset voidaan toteuttaa loppuosataulukolla.

### Hahmonsovitus

Tarkka hahmonsovitus tehdään tehostetulla binäärihaulla.

- Se on lähes yhtä nopea ( $\mathcal{O}(m + \log n + occ)$ ) kuin hahmonsovitus loppuosapuussa ( $\mathcal{O}(m + occ)$ ).
- Toisin kuin loppuosapuussa, hakuaika pätee, vaikka aakkoston koko ei ole vakio.
- Esiintymien lukumäärä voidaan laskea ajassa  $\mathcal{O}(m + \log n)$ , sillä binäärihaku tuottaa osavälin, jonka koko on esiintymien lukumäärä. Loppuosapuussa tilanne on hieman monimutkaisempi (HT).

- On olemassa myös aivan toisenlainen hakumenetelmä, takaperinhaku, jolla myös loppuosataulukossa päästään hakuaikaan  $\mathcal{O}(m + occ)$  vakioaakkostolla.

### Merkkijonon sisäinen rakenne

LCP-arvoista saadaan tietoa jonon rakenteesta.

- Pisimmän toistuvan osajonon pituus on LCP-taulukon maksimi.
- Erilaisten osajonojen lukumäärä saadaan kaavasta.

$$\frac{n(n+1)}{2} + 1 - \sum_{i=0}^{n-1} LCP[i]$$

### Yleistetty loppuosataulukko

- Useamman merkkijonon yleistetty loppuosataulukko voidaan määritellä kuten yleistetty loppuosapuun ja sen sovellukset ovat paljolti samat.
- Esimerkiksi kahden jonon  $S$  ja  $T$  pisin yhteinen osajono löydetään etsimällä jonon  $R = S\mathcal{L}T\$$  loppuosataulukosta suurin  $LCP[i]$ , jolla  $R_{SA[i]}$  ja  $R_{SA[i+1]}$  ovat eri jonoista.

### Vahvistettu loppuosataulukko

M. I. Abouelhoda & S. Kurtz & E. Ohlebusch:  
Replacing Suffix Trees with Enhanced Suffix Arrays. *J. Discrete Algorithms*, 2(1):53–86, 2004.

- Vahvistettu loppuosataulukko lisää normaaliin loppuosataulukkoon ja LCP-taulukkoon vielä pari lisätaulukkoa, jotka tekevät siitä käytännössä loppuosapuun.
- Merkittävin lisä ovat loppuosalinkit, jotka normaalista loppuosataulukosta puuttuvat.

### RMQ-esikäsittely

- Taulukon *osaväliminimi* (range minimum) on pienin arvo annetulla taulukon osavälillä. Mikä tahansa taulukko voidaan esikäsitellä lineaarisessa ajassa siten, että jokaiseen osavälikyselyyn voidaan vastata vakioajassa.
- LCP-taulukossa tehtävä osavälikysely (RMQ, range minimum query) antaa kahden loppuosan pisimmän yhteisen alkuosan.

**Lemma 4.6.1** *Kahden loppuosan  $S_i < S_j$  pisimmän yhteisen alkuosan pituus  $lcp(S_i, S_j)$  on  $\min\{LCP[k] \mid SA^{-1}[i] \leq k < SA^{-1}[j]\}$ .*

*Todistus.*

- Olkoon  $i' = SA^{-1}[i]$  ja  $j' = SA^{-1}[j]$ .
- Siis  $S_i = S_{SA[i']} < S_{SA[i'+1]} < \dots < S_{SA[j'-1]} < S_{SA[j']} = S_j$ .
- Väite seuraa nyt lemmän 3.7.2 yleistyksestä (induktiolla, HT).

□

- RMQ-esikäsitellyllä loppuosataulukolla voidaan korvata LCA-esikäsitelty loppuosapuu.
- RMQ-esikäsitely on yksinkertaisempi kuin LCA-esikäsitely.

### Burrows–Wheeler-muunnos

M. Burrows & D. Wheeler: A Block Sorting Lossless Data Compression Algorithm. Technical Report 124, Digital Equipment Corporation, 1994.

- Yksi tärkeimmistä tiedontiivistyksen menetelmistä on Burrows–Wheeler-muunnos. Sillä on läheinen yhteys loppuosataulukoon.
- Tavallisin tapa määritellä merkkijonon  $S$  Burrows–Wheeler muunnos on  $S$ :n kaikkien rotaatioiden akkosjärjestetyn matriisin viimeisenä sarakkeena. Se voidaan myös määritellä loppuosataulukon avulla.

- Olkoon  $S[0, n)$  merkkijono. Oletetaan lisäksi, että  $S[n] = \$$ , missä  $\$$  on kaikkia merkkejä pienempi. BW-muunnos  $BWT[0, n]$  on

$$BWT[i] = \begin{cases} \$ & \text{jos } SA[i] = 0 \\ S[SA[i] - 1] & \text{muuten} \end{cases}$$

**Esimerkki 4.6.2** Jonon *banana\$* BW-muunnos on *annb\$aa*.

	SA
\$banan a	\$
a\$bana n	a\$
ana\$ba n	ana\$
anana\$ b	anana\$
banana \$	banana\$
na\$ban a	na\$
nana\$b a	nana\$

- BW-muunnos on kääntyvä, eli alkuperäinen teksti voidaan palauttaa muunnoksesta.

BW-muunnosta käytetään tiedontiivistyksessä.

- Muunnettu teksti on helpommin tiivistyvä kuin alkuperäinen, koska se sisältää paljon paikallista toistoa.
  - Esim. englanninkielisessä tekstissä **he**-alkuiset rotaatiot tulevat järjestämisen jälkeen peräkkäin, jolloin viimeisessä sarakkeessa on vastaavalla kohdalla paljon kirjaimia **t** ja **T**.
- BW-muunnokseen perustuvat myös useat tiivistetyt loppuosarakenteet.