

# Sufixové polia

kuko

24.11.2018

Vybrané partie z dátových štruktúr

## Suffix array

Array of suffixes in lexicographic order

(assume  $\$ < a \quad \forall a \in \Sigma$ )

i	0	1	2	3	4	5	6
S[i]	b	a	n	a	n	a	\$

i	0	1	2	3	4
S[i]	a	a	a	a	\$

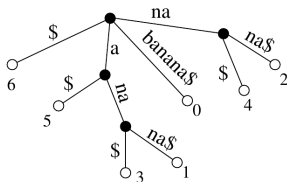
i	SA[i]	Suffix
0	<b>6</b>	\$
1	<b>5</b>	a\$
2	<b>3</b>	ana\$
3	<b>1</b>	anana\$
4	<b>0</b>	banana\$
5	<b>4</b>	na\$
6	<b>2</b>	nana\$

i	SA[i]	Suffix
0	<b>4</b>	\$
1	<b>3</b>	a\$
2	<b>2</b>	aa\$
3	<b>1</b>	aaa\$
4	<b>0</b>	aaaa\$

## From suffix array to suffix tree

T = banana\$

i	SA[i]	L[i]	suffix
0	<b>6</b>	0	\$
1	<b>5</b>	1	a\$
2	<b>3</b>	3	ana\$
3	<b>1</b>	0	anana\$
4	<b>0</b>	0	banana\$
5	<b>4</b>	2	na\$
6	<b>2</b>	-	nana\$



## Longest common prefix

- $\text{lcp}(A, B)$  = the length of longest common prefix of strings  $A$  and  $B$
- $\text{LCP}(i, j) = \text{lcp}(T[\text{SA}[i]..n], T[\text{SA}[j]..n])$   
i.e. lcp of two suffixes in a suffix array

$i$	$\text{SA}[i]$	Suffix
0	6	\$
1	5	a\$
2	3	ana\$
3	1	anana\$
4	0	banana\$
5	4	na\$
6	2	nana\$

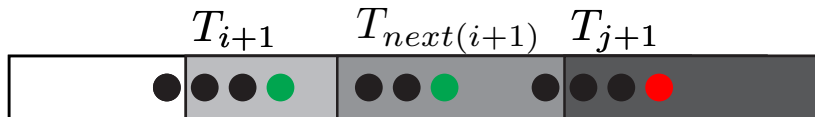
$$\text{LCP}(2,3) = \text{lcp}(\text{ana}\$, \text{anana}\$) = 3$$

$$\text{LCP}(2,5) = \text{lcp}(\text{ana}\$, \text{na}\$) = 0$$

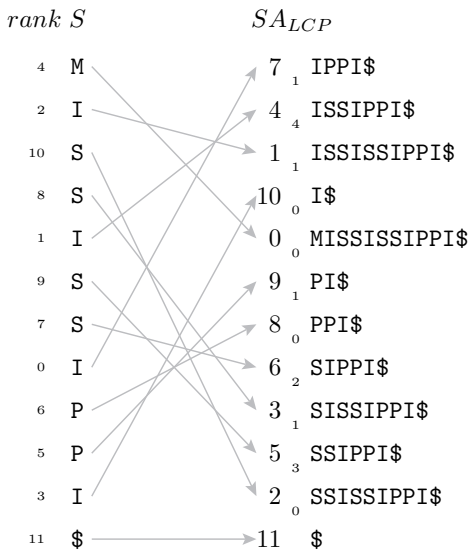
```
1 h = 0;
2 for (i=0; i <=n; i++) {
3     if (rank[i] < n) {
4         j = next(i);    // SA[rank[i] + 1];
5         // compare suffixes S[i..] and S[j..]
6         h = 0;
7         while (S[i+h] == S[j+h]) ++h;
8         // we have found first mismatch
9         L[i] = h;
10    }
11 }
```











```
1 h = 0;
2 for (i=0; i <=n; i++) {
3     if (rank[i] < n) {
4         j = next(i);    // SA[rank[i] + 1];
5         // compare suffixes S[i..] and S[j..]
6         h = 0;
7         while (S[i+h] == S[j+h]) ++h;
8         // we have found first mismatch
9         L[i] = h;
10    }
11 }
```

```
1 h = 0;
2 for (i=0; i <=n; i++) {
3     if (rank[i] < n) {
4         j = next(i);    // SA[rank[i] + 1];
5         // compare suffixes S[i..] and S[j..]
6         // assume they have >= h characters in common
7         while (S[i+h] == S[j+h]) ++h;
8         // we have found first mismatch
9         L[i] = h;
10        if (h > 0) --h;
11    }
12 }
```

## Inverse of a suffix array

Array rank such that  $\text{rank}[i] = x \iff \text{SA}[x] = i$

Can be computed in  $O(n)$  from SA:

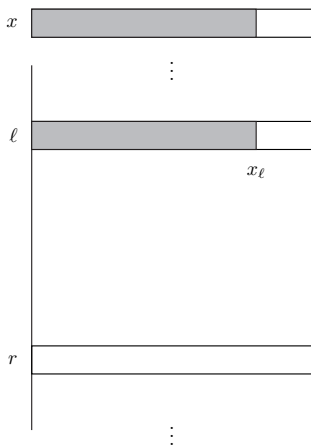
```
1 for (i = 0; i <= n; i++) {  
2   rank[SA[i]] = i;  
3 }
```

Go from suffix to its position in the suffix array, its neighbors, etc.

## Vyhľadávanie

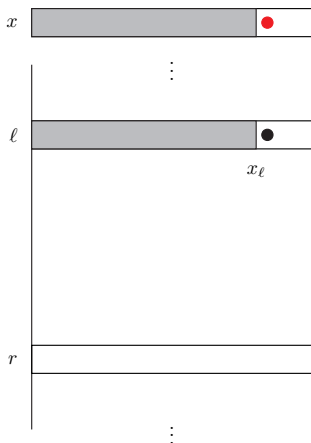
- jednoduché binárne vyhľadávanie:  $O(m \log n)$
- ak poznáme  $\text{lcp}(s_i, s_j)$  pre ľubovoľné dva sufiky, vieme zrýchliť na  $O(m + \log n)$

INVARIANTY:  $\ell < x \leq r$ ,  $x_\ell = \text{lcp}(x, \ell)$ ,  $x_r = \text{lcp}(x, r)$



Obr.: Potiaľto sa  $x$  a  $\ell$  rovnajú.

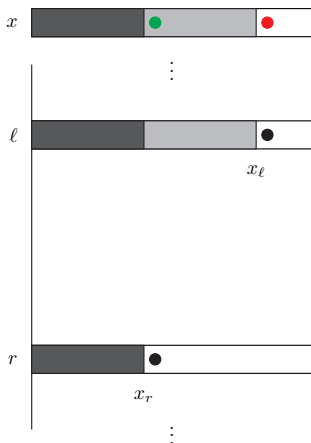
INVARIANTY:  $\ell < x \leq r$ ,  $x_\ell = \text{lcp}(x, \ell)$ ,  $x_r = \text{lcp}(x, r)$



Obr.:  $\ell[x_\ell] < x[x_\ell]$ .

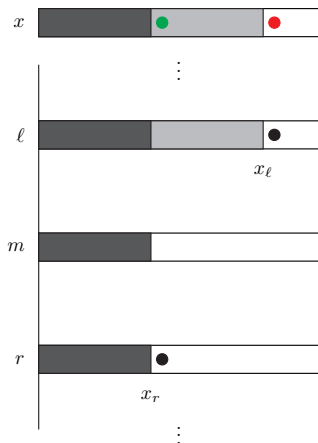


INVARIANTY:  $\ell < x \leq r$ ,  $x_\ell = \text{lcp}(x, \ell)$ ,  $x_r = \text{lcp}(x, r)$



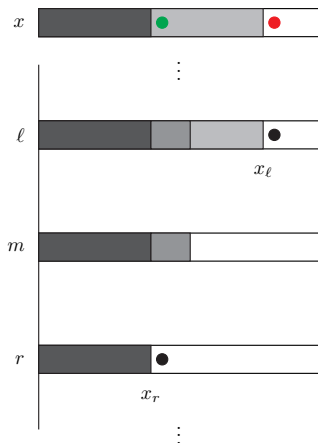
Obr.: Potiaľto sa  $x$  a  $r$  rovnajú. (Predpokladajme  $x_\ell \geq x_r$ )

INVARIANTY:  $\ell < x \leq r$ ,  $x_\ell = \text{lcp}(x, \ell)$ ,  $x_r = \text{lcp}(x, r)$



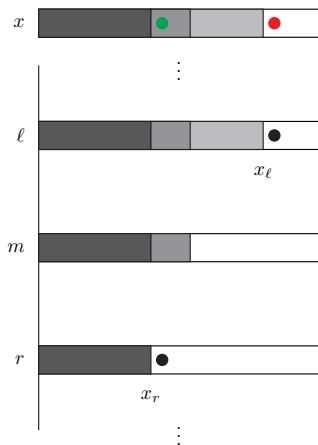
Obr.:  $r[x_r] > x[x_r]$ .

INVARIANTY:  $\ell < x \leq r$ ,  $x_\ell = \text{lcp}(x, \ell)$ ,  $x_r = \text{lcp}(x, r)$



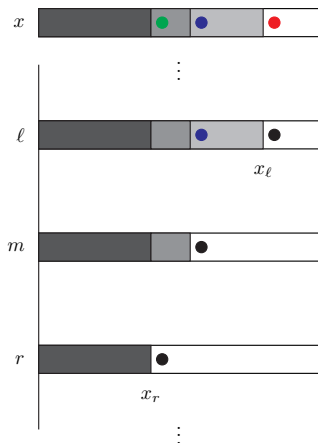
Obr.: Pozrime sa na prostredný sufix; aký je  $p = \text{lcp}(\ell, m)$ ?

INVARIANTY:  $\ell < x \leq r$ ,  $x_\ell = \text{lcp}(x, \ell)$ ,  $x_r = \text{lcp}(x, r)$



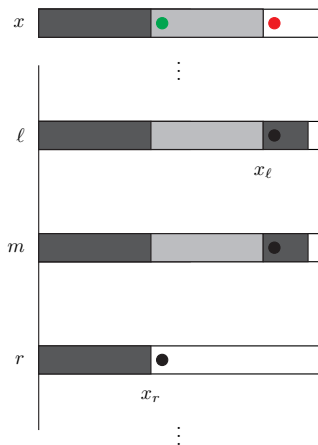
Obr.: Prípád 1: Ak  $p < x_\ell \dots$

INVARIANTY:  $\ell < x \leq r$ ,  $x_\ell = \text{lcp}(x, \ell)$ ,  $x_r = \text{lcp}(x, r)$



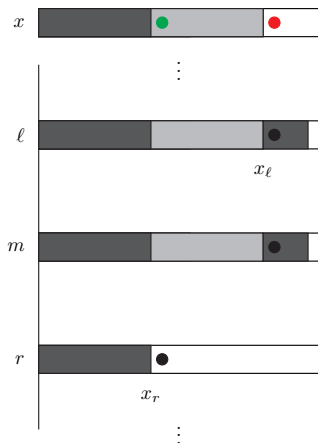
Obr.: Prípád 1: ... tak  $m[p] > \ell[p] = x[p]$ , teda  $m > x$ .

INVARIANTY:  $\ell < x \leq r$ ,  $x_\ell = \text{lcp}(x, \ell)$ ,  $x_r = \text{lcp}(x, r)$



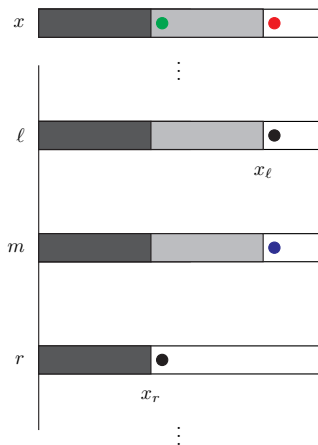
Obr.: Prípád 2: Ak  $p > x_\ell, \dots$

INVARIANTY:  $\ell < x \leq r$ ,  $x_\ell = \text{lcp}(x, \ell)$ ,  $x_r = \text{lcp}(x, r)$



Obr.: Prípád 2: ... tak  $m[x_\ell] = \ell[x_\ell] < x[x_\ell]$ , teda  $m < x$ .

INVARIANTY:  $\ell < x \leq r$ ,  $x_\ell = \text{lcp}(x, \ell)$ ,  $x_r = \text{lcp}(x, r)$



Obr.: Prípád 3: Ak  $p = x_\ell$ , začneme porovnávať písmená.



Invarianty:

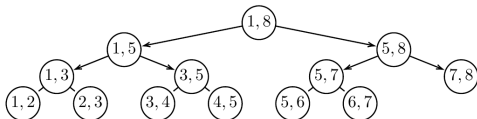
- $l < x \leq r$
- $x_l = \text{lcp}(x, l)$
- $x_r = \text{lcp}(x, r)$

Algoritmus:

- Ak  $x_l \geq x_r$ :
  - nech  $m$  je prostredný sufix a  $p = \text{lcp}(l, m)$
  - ak  $p < x_l$ :  $r \leftarrow m$ ,  $x_r \leftarrow p$
  - ak  $p > x_l$ :  $l \leftarrow m$
  - ak  $p = x_l$ :
    - začneme porovnávať písmená  $x$  a  $m$  od pozície  $x_l$
    - podľa výsledku nastavíme  $l, x_l$  alebo  $r, x_r$
- Ak  $x_l \leq x_r$  – symetricky

### LCP values for algorithm 3

Which values are needed?  $LCP(L, k)$  or  $LCP(R, k)$



$2n - 1$  LCP values needed

Let  $L[i] = LCP(i, i + 1)$ , precompute to an array in  $O(n)$  (later)

For  $j - i > 1$ :

$$\begin{aligned} LCP(i, j) &= \min\{LCP(k, k + 1) \mid k = i \dots j - 1\} \\ &= \min\{LCP(i, x), LCP(x, j)\} \text{ for any } x \in \{i + 1, \dots, j - 1\} \end{aligned}$$

## Konštrukcia SA

- qsort –  $O(n^2 \log n)$
- radix-sort / trie –  $O(n^2)$
- lepšie?

## Myšlienka:

- sufix sufixu je sufix
- ak máme pole utriedené podľa prvých  $K$  písmen, vieme ho ľahko zotriediť podľa prvých  $2K$  písmen
- budeme mať  $\log n$  fáz; v  $k$ -tej fáze triedime podľa prvých  $2^k$  písmen

Myšlienka:

- sufix sufixu je sufix
- ak máme pole utriedené podľa prvých  $K$  písmen, vieme ho ľahko zotriediť podľa prvých  $2K$  písmen
- budeme mať  $\log n$  fáz; v  $k$ -tej fáze triedime podľa prvých  $2^k$  písmen

# Manberov algoritmus

	index sort	inverse
0 babaaaabcbabaaaaa0	17 0	0 14
1 abaaaabcbabaaaaa0	16 a0	1 9
2 baaaabcbabaaaaa0	15 aa0	2 12
3 aaaabcbabaaaaa0	14 aaa0	3 4
4 aaabcbabaaaaa0	3 aaaabcbabaaaaa0	4 7
5 aabcbabaaaaa0	12 aaaaa0	5 8
6 abcbabaaaaa0	13 <b>aaaa0</b>	6 11
7 bcbabaaaaa0	4 <b>aaabcbabaaaaa0</b>	7 16
8 cbabaaaaa0	5 aabcbabaaaaa0	8 17
9 babaaaaa0	1 abaaaabcbabaaaaa0	9 15
10 abaaaaa0	10 abaaaaa0	10 10
11 baaaaa0	6 abcbabaaaaa0	11 13
12 aaaaa0	2 baaaabcbabaaaaa0	12 5
13 aaaa0	11 baaaaa0	13 6
14 aaa0	0 babaaa <b>abcbabaaaaa0</b>	14 3
15 aa0	9 bab <b>aaaaa0</b>	15 2
16 a0	7 bcbabaaaaa0	16 1
17 0	8 cbabaaaaa0	17 0

$0 + 4 = 4$  →

$9 + 4 = 13$  →

13 < 4 (because 6 < 7) so 9 < 0

Obr.: Triedime podľa prvých 1,2,4,8,16,32,... písmen.

- nech  $\text{rank}[i] = j$ , ak je sufix  $s_i$  v abecednom poradí  $j$ -ty (podľa prvých  $2^k$  písmen; ak majú dva sufíxy rovnakých prvých  $2^k$  písmen, ranky budú rovnaké)
- 1 fáza:
  - zotriedime trojice  $(\text{rank}[i], \text{rank}[i + 2^k], i)$



- qsort –  $O(n^2 \log n)$
- radix-sort / trie –  $O(n^2)$
- Manberov algoritmus –  $O(n \log n)$
- lepšie?

1 2 4 5 7 8 10 11 13  
A B B A B A B A B B A \$ \$ \$

0 3 6 9 12  
A B B A B A B A B B A \$ \$ \$

0.  $S_{10} = A$$$$
1.  $S_5 = ABABBA$$$$
2.  $S_7 = ABBA$$$$
3.  $S_2 = BABABABBA$$$$
4.  $S_4 = BABABBA$$$$
5.  $S_8 = BBA$$$$
6.  $S_1 = BBABABABBA$$$$

Obr.: Rozdelíme na sufíxy na pozíciách nedeliteľných vs. deliteľných 3.

1 2 4 5 7 8 10 11 13  
A B B A B A B A B B A \$ \$ \$

0 3 6 9 12  
A B B A B A B A B B A \$ \$ \$

0.  $S_{10}$
1.  $S_5$
2.  $S_7$
3.  $S_2$
4.  $S_4$
5.  $S_8$
6.  $S_1$

**Obr.:** Rekurzívne utriedime pozície  $\equiv 1, 2 \pmod{3}$ .  
Čo pozície  $\equiv 0 \pmod{3}$ ?

1 2 4 5 7 8 10 11 13  
A B B A B A B A B B A \$ \$ \$

0 3 6 9 12  
A B B A B A B A B B A \$ \$ \$

0.  $S_{10}$
1.  $S_5$
2.  $S_7$
3.  $S_2$
4.  $S_4$
5.  $S_8$
6.  $S_1$

0.  $S_3 = AS_4$
1.  $S_0 = AS_1$
2.  $S_9 = BS_{10}$
3.  $S_6 = BS_7$

**Obr.:** Odrolujeme jedno písmeno; zvyšok je poz.  $\equiv 1 \pmod{3}$ .

1 2 4 5 7 8 10 11 13  
A B B A B A B A B B A \$ \$ \$

0 3 6 9 12  
A B B A B A B A B B A \$ \$ \$

0.  $S_{10}$
1.  $S_5$
2.  $S_7$
3.  $S_2$
4.  $S_4$
5.  $S_8$
6.  $S_1$

0.  $S_3$
1.  $S_0$
2.  $S_9$
3.  $S_6$

Obr.: Ako tieto dve utriedené polia zmerge-ujeme?

1 2 4 5 7 8 10 11 13  
A B B A B A B A B B A \$ \$ \$

0 3 6 9 12  
A B B A B A B A B B A \$ \$ \$

0.  $S_{10} = AS_{11}$

1.  $S_5 = ABS_7$

2.  $S_7 = AS_8$

3.  $S_2 = BAS_4$

4.  $S_4 = BS_5$

5.  $S_8 = BBS_{10}$

6.  $S_1 = BS_2$

0.  $S_3 = AS_4 = ABS_5$

1.  $S_0 = AS_1 = ABS_2$

2.  $S_9 = BS_{10} = BAS_{11}$

3.  $S_6 = BS_7 = BAS_8$

Obr.: Odrolujeme jedno alebo dve písmená.

Zložitosť:

- $T(n) =$ 
  - $T(\frac{2}{3}n)$  – rekurzívne volanie
  - $+O(n)$  – triedenie  $\equiv 0 \pmod{3}$ 
    - $+O(n)$  – merge

Moment: ale aké rekurzívne volanie? Je to tá istá úloha?

Zložitosť:

- $T(n) =$ 
  - $T(\frac{2}{3}n)$  – rekurzívne volanie
  - $+O(n)$  – triedenie  $\equiv 0 \pmod{3}$ 
    - $+O(n)$  – merge

Moment: ale aké rekurzívne volanie? Je to tá istá úloha?



A B B A B A B A B B A \$ \$ \$

Obr.: Začneme s pôvodným stringom.

1                    4                    7                    10                    13  
A B B A B A B A B B A \$ \$ \$

1                    4                    7                    10  
B B A B A B A B B A \$ \$

Obr.: Vezmeme jednu kópiu od pozície 1...

1 2 4 5 7 8 10 11 13  
A B B A B A B A B B A \$ \$ \$

1 4 7 10 2 5 8  
B B A B A B A B B A \$ \$ B A B A B A B B A

Obr.: ... a pridáme kópiu od pozície 2...

Ak budeme každú trojicu znakov považovať za 1 písmeno, tak sufíxy tohto stringu zodpovedajú sufíxom na poz.  $\equiv 0, 1 \pmod{3}$  v pôvodnom stringu.

Problém: Veľká abeceda – bude porovnávanie znakov v  $O(1)$ ?

1 2 4 5 7 8 10 11 13  
A B B A B A B A B B A \$ \$ \$

1 4 7 10 2 5 8  
B B A B A B A B B A B A B B A

Obr.: ... a pridáme kópiu od pozície 2...

Ak budeme každú trojicu znakov považovať za 1 písmeno, tak sufíxy tohto stringu zodpovedajú sufíxom na poz.  $\equiv 0, 1 \pmod{3}$  v pôvodnom stringu.

Problém: Veľká abeceda – bude porovnávanie znakov v  $O(1)$ ?

1 2 4 5 7 8 10 11 13  
A B B A B A B A B B A \$ \$ \$

1 4 7 10 2 5 8  
B B A B A B A B B A B A B B A

Obr.: ... a pridáme kópiu od pozície 2...

Ak budeme každú trojicu znakov považovať za 1 písmeno, tak sufíxy tohto stringu zodpovedajú sufíxom na poz.  $\equiv 0, 1 \pmod{3}$  v pôvodnom stringu.

Problém: Veľká abeceda – bude porovnávanie znakov v  $O(1)$ ?

1 2 4 5 7 8 10 11 13  
A B B A B A B A B B A \$ \$ \$

A\$\$ → a

ABA → b

ABB → c

BAB → d

BBA → e

1 4 7 10 2 5 8  
B B A B A B A B B A \$ \$ B A B A B A B B A

Obr.: Riešenie: trojíc je málo, takže si ich môžeme prečíslovať.

1 2 4 5 7 8 10 11 13  
A B B A B A B A B B A \$ \$ \$

A\$\$ → a

ABA → b

ABB → c

BAB → d

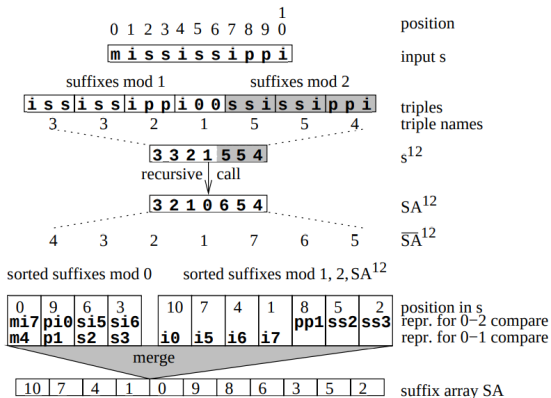
BBA → e

1 4 7 10 2 5 8  
B B A B A B A B A B B A

e d c a d b e

**Obr.:** Takto ostane abeceda malá. Rekurzívne sa zavoláme na string edcadbe. Výsledné sufixové pole musíme následne prečíslovať: sufixy 0,1,2,3 zodpovedajú pozíciám 1,4,7,10 a sufixy za polovicou: 4,5,6 zodpovedajú pozíciám 2,5,8.

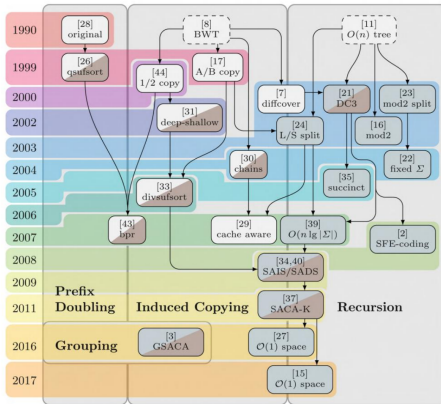
# Konštrukcia SA



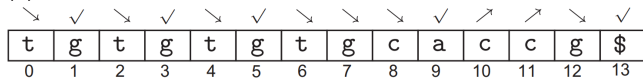


Algoritmus:

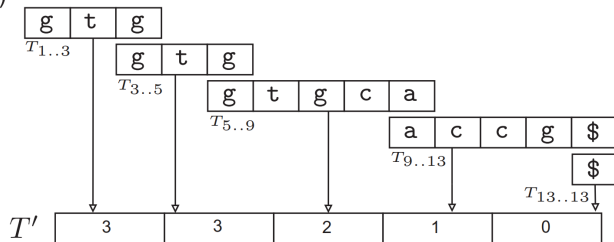
- $SA^{1,2}$  – zotriedime pozície  $\equiv 1, 2 \pmod{3}$ :
  - vytvoríme string  $s'$  dĺžky  $n' = \frac{2}{3}n$ :
    - vytvoríme  $s^{1,2}$  tak, že zretazíme  $s[1\dots]$  a  $s[2\dots]$
    - zredukujeme abecedu: radix-sortom utriedime všetky trojice znakov v  $s^{1,2}$  a prepíšeme ich na znaky  $1, \dots, k$
  - rekurzívne zavoláme  $SA(s')$
  - vo výsledku prečíslujeme indexy:
    - ak  $i < \lceil n'/2 \rceil \rightarrow 3i + 1$
    - ak  $i \geq \lceil n'/2 \rceil \rightarrow 3(i - \lceil n'/2 \rceil) + 2$
- $SA^0$  – zotriedime pozície  $\equiv 0 \pmod{3}$ :
  - radix-sortom podľa prvého písmena a pozície v  $SA^{1,2}$
- zmergujeme  $SA^{1,2}$  a  $SA^0$ 
  - podľa prvého alebo prvých dvoch písmen a pozície v  $SA^{1,2}$

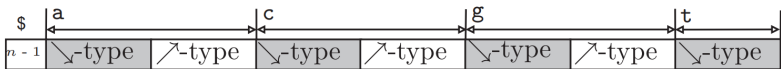


(a)



(b)





	a	c			g					t				
13	9	-1	-1	-1	-1	-1	5	3	1	-1	-1	-1	-1	





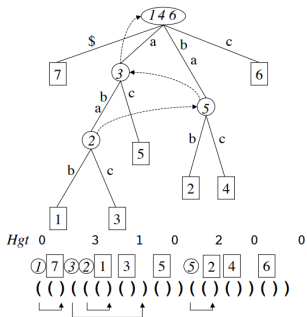
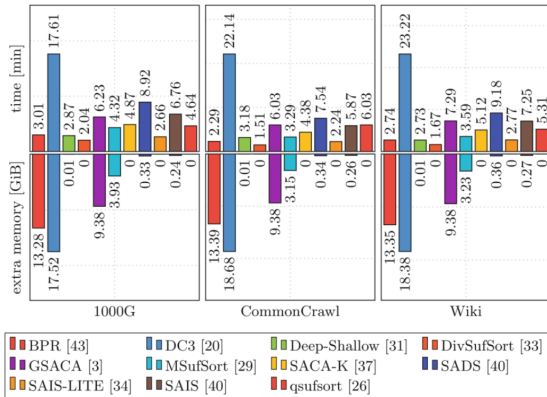
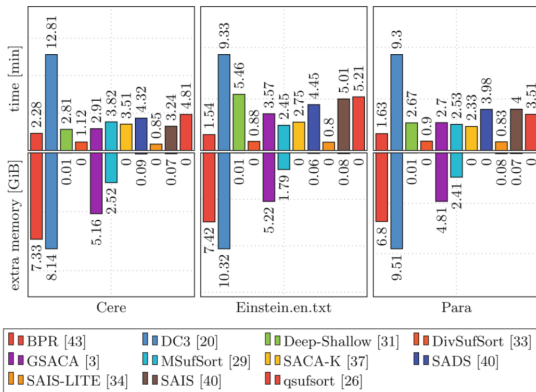
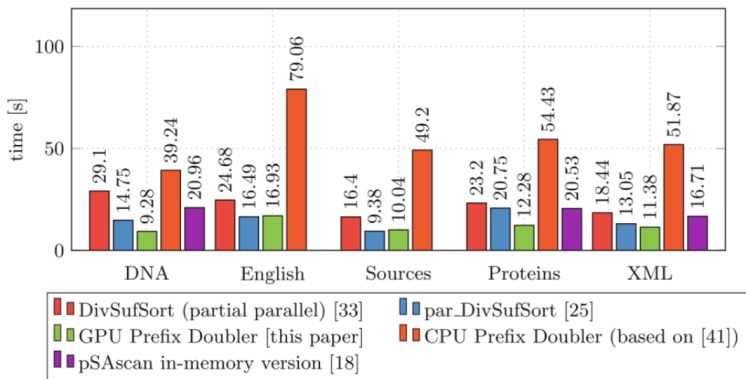


Figure 1: The suffix tree for “ababac\$,” and its balanced parentheses representation. Suffix links are shown in dotted line.









$i$	1	2	3	4	5	6	7
$Hgt$	0	3	1	0	2	0	0
$SA$	7	1	3	5	2	4	6
$SA + Hgt$	7	4	4	5	4	4	6

Figure 2: How to create a sorted sequence from  $Hgt$ .