

Sufixový pole

Ds pro textová alq.: zpracování textů, problém matky, binární strom, ...

vyhledávání podřetězec B , $m = |B|$ v řetězci A , $n = |A|$:

- tabul - karp (volující károu) ... $O(n \cdot m)$ čas
 - kuchař - Morris - Pratt
 - Aho - Corasick
- $O(n + m)$ } vhodné pro pevné B

Co když A je pevné? Uvidíme: lze $O(\log n + m)$ (přesně #y'syho na výstup)

Zusatz

$d = d[0] \dots d[n-1]$ řetězec

$d[i:j] = d[i] \dots d[j-1]$ podřetězec délky $j-i$, pro j i podřetězí řetězce \in

$d[0:j]$... prefix délky j

$d[i:n]$... suffix začínající $d[i]$ délky $n-i$

$d \leq B$... lex. uspořádání: d prefix B nebo existuje i tiz

lze vyhledat pomocí $d[i:j] = B[i:j]$ a $d[i] < B[i]$
a $i, B \neq$

Př.

$d = \text{bananas}, n = 7$

0 1 2 3 4 5 6 7

i	suffix	$S[i:]$	$R[i]$	$L[i]$
0	ϵ	7	4	0
1	ananas	1	1	3
2	anas	3	5	1
3	as	5	2	0
4	bananas	0	6	0
5	nanas	2	3	2
6	nas	4	7	0
7	s	6	0	-

číslo: pro lib. S , vyšetřte S v daném intervalu v S

\Rightarrow pomocí bin. vyhledávání lze najít všechny vyšetřte S v čase $O(m \log n + p)$

vyšetřte

Def: Suffixové pole pro d je permutace $S[0..n]$ t.j. $d[S[i]:] < d[S[i+1]:]$

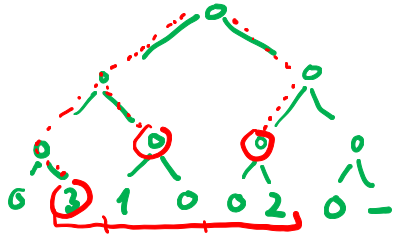
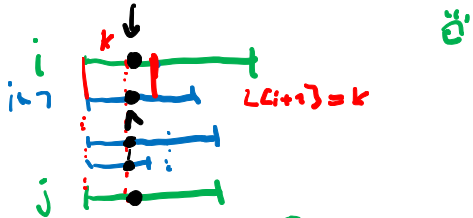
je také $0 \leq i < n$, t.j. $S[i]$ je začátek intervalu lex. nejmenšího suffixu d .

• Rankové pole $R[0..n]$ je inverzní permutace k S , t.j. $R[i] = \text{lokace}$ lex. nejmenší suffix je $d[i:]$ (zjistit v $O(n)$ ze S)

• LCP pole $L[0..n-1]$, $L[i] = \text{LCP}(d[S[i]:], d[S[i+1]:])$, t.j. délka největšího společného prefixu

$LCP(p, q) = \max k \text{ s.t. } p[1:k] = q[1:k]$

LCP prin lib. din suffix



$$LCP(d[S[i]:], d[S[j]:]) = \min(L[i], \dots, L[j-1])$$

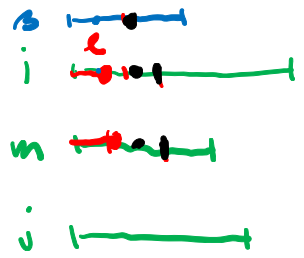
LCP $\geq k$ traverse, LCP $\leq k$ urbat - postapunct
 pisun va pade: k+1 si neblanaje, urtiti naste

stare prin intervalul dat

$$\Rightarrow \text{caz } O(\log n) \text{ (dat)}$$

$$\text{BUILD } O(n)$$

Idea: Zădări LCP prin lib. din suffix, la unizat lex. nejm. suffix d, k
 un prefix B si caz $O(m + \log n)$.



gătim interval $[i, j]$, $l = LCP(B, d[S[i]:])$

$$m = \lfloor (i+j)/2 \rfloor$$

if $l > LCP(d[S[i]:], d[S[m]:])$ continue $[i, m]$

$<$
 $=$ parcurgem B a $d[S[m]:]$ ad parca l
 (l si p'iment)

\Rightarrow všechny výskyt \exists v d lze načíst potom v čase $O(m + \log n)$.

Dokonce, pro lib. m $LCP(i, m)$ a $LCP(m, j)$ potřeby se pro nejvyšší
 jedno i a nejvyšší jedno j . ("LCP-LR" pole lze v čase $O(n)$)

Aplikační LCP pole

i) histogram pro k -gramy = # výskytů pro každý podřetězec délky k
 v čase $O(n)$ (k -gram)

Př. $d = \text{bananas}$, $n = 7$
 $0\ 1\ 2\ 3\ 4\ 5\ 6\ 7$

i	suffix	$S[i]$	$R[i]$	$L[i]$
0	ϵ	7	4	0
1	[anas	1	1	3
2	[anas	3	5	1
3	[as	5	2	0
4	bananas	0	6	0
5	[nanas	2	3	2
6	[nas	4	7	0
7	s	6	0	-

interval $[i, j]$ t.j.

$L[i-1] < k, L[i] \geq k$ & i < j
 $L[j] < k$

$k=2$

an 2

as 1

ba 1

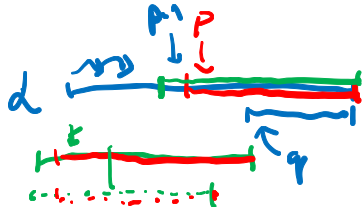
na 2

2) nejdelší opakující se podřetězec = $\max_i L[i]$ v čase $O(n)$

3) nejdelší společný podřetězec v $d, s = \max L[i]$
 $i, S[i] = S[i+1]$
 $2 \cdot d/3$
 LCP po $d \neq s$

Konstrukce LCP-pole (kasai)

Nějme $S, R \rightarrow L$ v čase $O(n)$

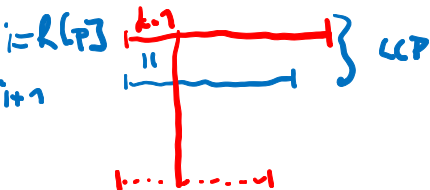


if $k = L[R[p-1]]$,
 pak $L[R[p]] \geq k-1$

BUILD LCP

```

k := 0
for p = 0 ... n-1
  k ← max(0, k-1)
  i ← R[p]
  q ← S[i+1]
  while (p+k < n) ∧ (q+k < n)
    ∧ (S[p+k] = S[q+k])
  
```



if: min while ... $O(n)$ čas
 \neq přířekni 1 = # odčítání 1 + hodnoty konstantní
 \Rightarrow celkem $O(n)$

— počítání $\leq k$ $t \leftarrow k+1$
 hodnota $L[i] \leftarrow k$

konstante suf. pole $\sim O(n \log n)$

Def: řetězce ρ, σ : $\rho \stackrel{k}{\leq} \sigma$ pokud $\rho[1:k] = \sigma[1:k]$
 $\rho \leq_k \sigma$ $\rho[1:k] \leq \sigma[1:k]$
 \uparrow lex.

$R_k[i] = \#$ množství sufixů vzhledem \leq_k než $\rho[i]$

$\exists \rho[1:i] \leq_k \rho[1:j] \Leftrightarrow \rho[1:i] <_k \rho[1:j]$ nebo $(\rho[1:i] =_k \rho[1:j] \wedge \rho[1+i:k] \leq_k \rho[1+j:k])$

$R_{2k}[i] \leq R_{2k}[j] \Leftrightarrow R_k[i] < R_k[j]$ nebo $(R_k[i] = R_k[j] \wedge R_k[i+k] \leq R_k[j+k])$
 $\Leftrightarrow (R_k[i], R_k[i+k]) \leq_{\text{lex}} (R_k[j], R_k[j+k])$

Alg: $k=2^0$ S_1, R_k podle prvních písmen $O(n \log n)$

$k \rightarrow 2k$ bucket sort $(R_k[i], R_k[i+k]) \dots$ 2 písmeny $O(n)$

\Rightarrow celkem $O(n \log n)$