

Práce s řetězci

$T \in \Sigma^m$... text

$P \in \Sigma^n$... hledané slovo

úkol: najdi všechny výskyty P v T .

Klasické řešení: Aho-Corasick - z P sestrojím konečný automat, který vyhledá výskyty P v T .
Čas $O(m+n)$.

typicky: T se nemění (fixní databáze) a přijde
řada dotazů P_1, P_2, \dots & $n \ll m$
→ chceme lepší řešení

suffixové pole

$T[1..m], T[2..m], T[3..m], \dots, T[m..m]$

... suffixy T

→ seřadím lexicograficky

→ $S[1..m]$

$S[i] = j$

$T[j..m]$ je lexicograficky
i-tý suffix.



↳ koncový symbol,
např. \0 u C

Pr: $T = \text{banana}\$$

$T[1..m] = \text{banana}\$$ $T[2..m] = \text{anan}\$$

$T[3..m] = \text{nan}\$$

⋮

$T[m..m] = \$$

$S[1] = \$$ $S[2] = \text{an}\$$ $S[3] = \text{anan}\$$ $S[4] = \text{banana}\$$

$S[5] = \text{n}\$$ $S[6] = \text{nan}\$$

vyhledání P pomocí suffixového pole

- binárně vyhledat P v suffixovém poli

$O(n \cdot \lg m)$

- lze zlepšit na $O(n + \lg m)$ pomocí Fw

pro najdení společný prefix řekněme $T[1..m]$
⋮
 $T[m..m]$

$\text{lcp}(a, b) = \max_k a[1..k] = b[1..k]$

- v čase $O(m \lg m)$ popř. $O(m)$ lze vybudovat
strukturu, která zodpovídá dotazy

$$LCP(i, j) = lcp(T[i..m], T[j..m])$$

v čase $O(1)$.

vybudování suffixového pole:

- lze vybudovat v čase $O(m \lg m)$ pomocí varianty bucket-sortu [Karp-Miller-Rosenberg '72]
- $\lg m$ fází, po fázi i jsou suffixy $T[j..m]$ seříděny podle prvního 2^i znaků.
- fáze $i=0$: seřídím $T[j..j]$, suffix $T[j..m]$
dostanu pořadí čísla $H[j]$
prvního výskytu znaku $T[j..j]$.
- fáze $i+1$: $T[j..m]$ označím dvojicí $(H[j], H[j+2^i])$
seřídím zůstane dvojice změníme
pořadí čísla $T[j..m]$ na pořadí
číslo dvojice $(H[j], H[j+2^i])$.
pořadí \rightarrow nové $H[j]$.
- po fázi $i = \lg m$ mi $H[j]$ udává pořadí čísla
 $T[j..m]$ mezi suffixy

\rightarrow rank $T[j..m]$.

$R[j] := H[j]$... pořadí čísla

$T[1..m]$ není fixy

$S[R[j]] = j$ tj. $R[..]$ je inverz $S[...]$.

- z R lze snadno doplnit S jedním přidáním
přes R : $\forall j=1..m: S[R[j]] := j$.

- i -tým fází algoritmu lze implementovat pomocí
bucket sortu / radix sortu v čase $O(m)$.

→ celkový čas $O(m \lg m)$ na spočítání $S[.]$.

LCP:

- pole $L[1..m]$; $L[i] = \text{lcp}(T[i], T[i+1])$.