

Problém s řetězi

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

$P \in \Sigma^n$... hledaná slova

úkol: najít všechny / všechny výskyty P v T

Klasika Aho - Corasick - sestrojím končící

automat $\rightarrow P$ a spuštím ho

na T

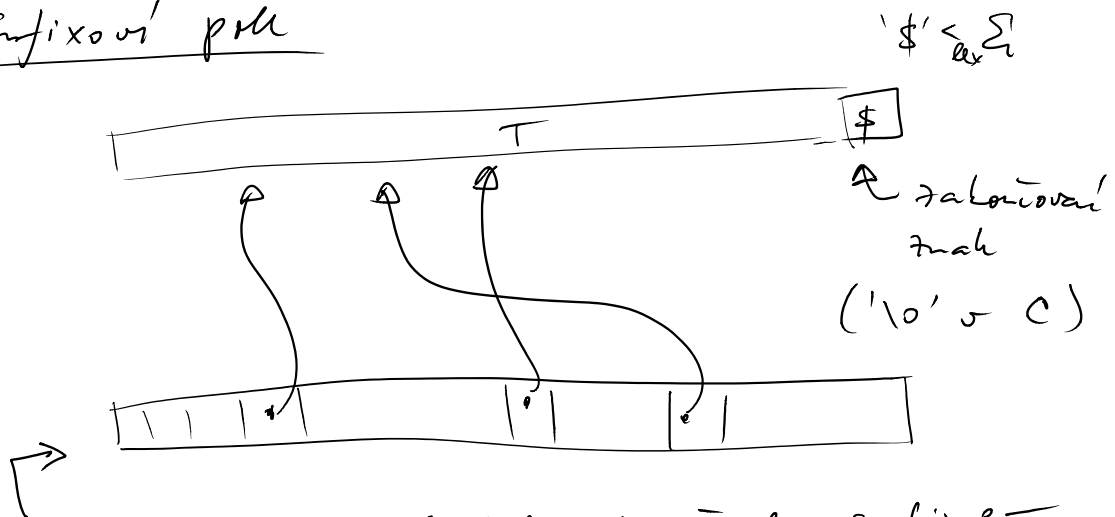
čas $\approx O(m+n)$.

chame lípe - typicky $m \gg n$ a T

je databáze, která se hemží

Řešení: sufixové pole, sufixový strom

Suffixové pole



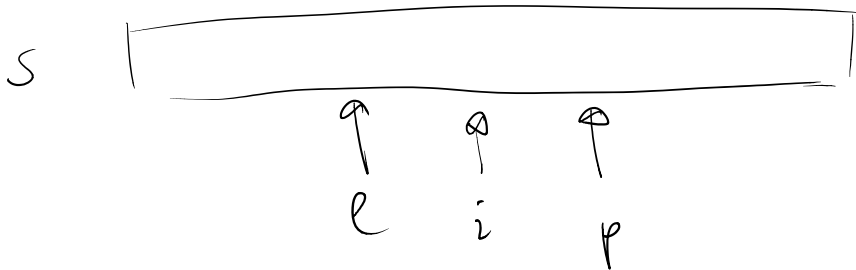
lexikograficky seřazené pole všech sufixů T ,
tj. pole řetězců $T[1..n], T[2..n], \dots, T[n..n]$

→ v poli lze binárně vyhledávat řetězec P

- hledám lexikograficky nejmenší řetězec v úložišti - roohn P.

• pokud se P shoduje s tímto řetězkem na pozici |P| pozicích, volal jsem výslyt P a T, bezprostředně následující řetězky též se shodující s P jsou další výslyt P.

binární vyhledávání - čas $O(n \cdot \log m)$



$l \leftarrow 1, p \leftarrow m$
dokud je $p > l + 1$

// hledá s přesností ±1

$i \leftarrow \lfloor \frac{l+p}{2} \rfloor$ (*)

pokud $S[i] <_{lex} P$ pak $l \leftarrow i$

jinak $p \leftarrow i$ ▽

return l

end.

... $\log m$ iterací, každé porovnání (*) stojí $O(m)$.

• lze zlepšit na $O(n + \log m)$

$u, v \in \Sigma^*$ $\text{lcp}(u, v)$ = délka největšího společného prefixu

předpokládáme, že známe zadarmo

$\text{lcp}(s(l), s(i))$ a $\text{lcp}(s(p), s(i))$

pro všechny trojice l, p, i (které

se mohou využít při binárním vyhledávání)

↳ tedy je práce $O(m)$, viz níže

• během binárního vyhledávání si udržuji:

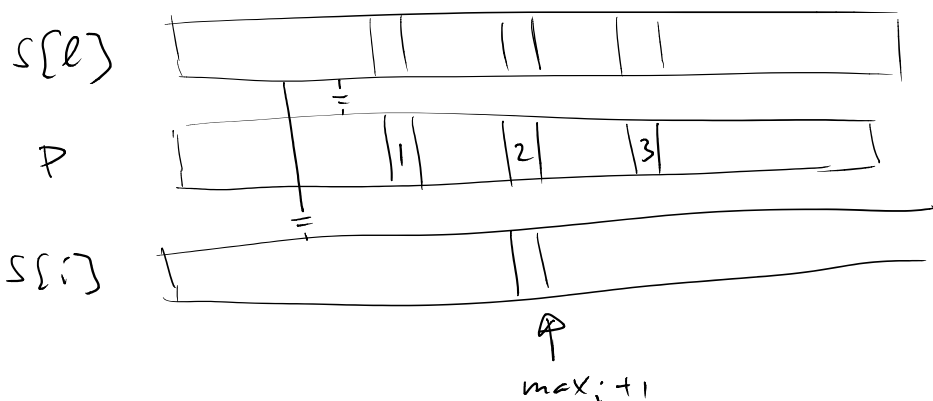
$$\max_l = \text{lcp}(s[l], P)$$

$$\max_p = \text{lcp}(s[p], P)$$

pokud $\max_l > \max_p$ porovnání (*) provedu

nahodově:

$$\max_i = \text{lcp}(s[l], s[i])$$



tři možnosti:

$$3) \max_l > \max_i \Rightarrow P <_{\text{lex}} s[i]$$

nebot' $S[l]$ a P se shodují
 až do (3) a $S[i]$ se od nich
 liší o $\max_i + 1$

$$\Rightarrow p \leftarrow i, \max_p \leftarrow \max_i$$

2) $\max_e = \max_i$: najdi první rozdílný znak
 $\max_i S[i]$ a P ze pozice $\max_i \rightarrow r$
 pokud $S[i][r] < P[r]$

$$\text{pok } l \leftarrow i, \max_e \leftarrow r - 1$$

$$\text{jinak } p \leftarrow i, \max_p \leftarrow r - 1$$

1) $\max_e < \max_i$: $l \leftarrow i, \max_e \leftarrow \max_i$

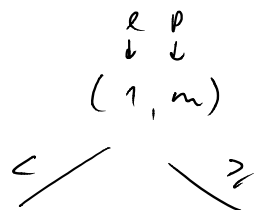
pokud $\max_e \leq \max_p$ postupuj: symetricky
 jako při: $\max_e > \max_p$.

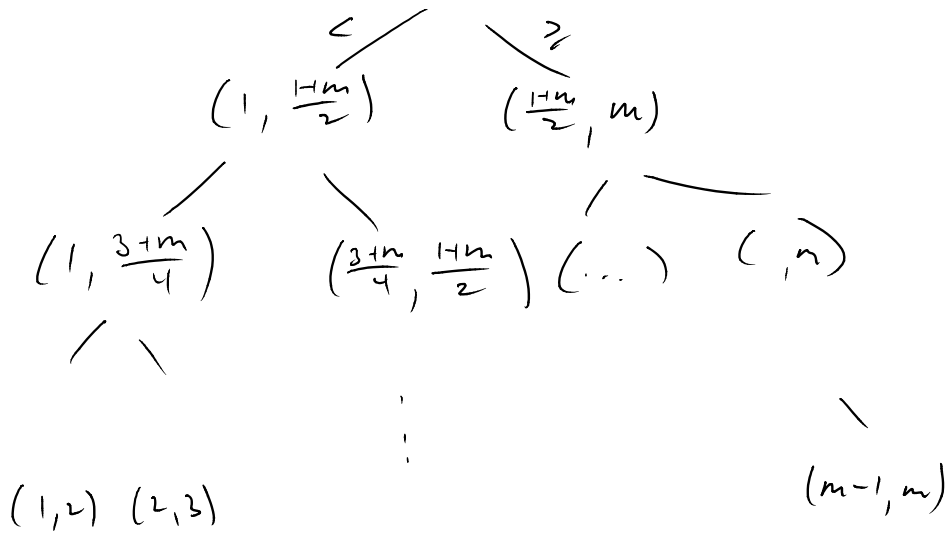
Pozn: pokud zjistíš, že P je prefix $S[i]$;
 uhlédáš tam zastavíš a navěd jsm
 ušlyt. Během binárního uhlédávání
 tak může předpokládat, že P není
 prefixem ani $S[l]$ ani $S[p]$.

čas na (*) je pak v směr $O(n + \log m)$,
 nebot' o P se posunují pouze doprava

• předpoklad, že zbudíme $lcp(S[l], S[i])$ a $lcp(S[p], S[i])$:

- ty si předpocítáš, je jich $O(m)$, tudíž
 je lze skladovat v rozumné paměti



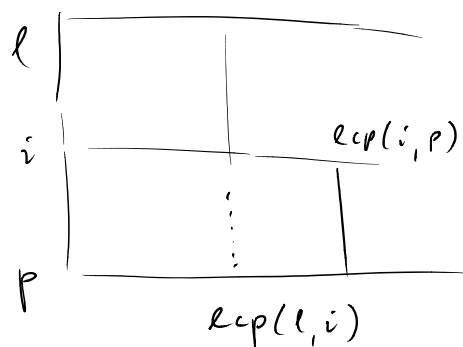


hranič l, p, i se posouvají jako při průchodu jedním z větví. Strom má $\leq m$ listů a $\leq m$ vnitřních uzlů $\Rightarrow \leq 2m$ kombinací l, p, i .

potřebné hodnoty lcp lze spočítat odspodu nahoru, pokud známe $lcp(i, i+1)$

$\forall i$. Platí totiž

$$lcp(l, p) = \min_{i \in [l, p]} lcp(i, i+1)$$



$$\Downarrow$$

$$lcp(l, p) = \min(lcp(l, i), lcp(i, p))$$

• $lcp(i, i+1) \quad \forall i = 1 \dots m-1$ lze spočítat

čas $O(m \lg m)$

• suffixový pole lze konstruovat v čase $O(m \cdot \lg m)$:

- varianta bucket-sortu [Karp-Miller-Rosenberg '72]

- $\lg m$ fází

• v multi-fázi rozdělím řetěz do bucketů podle prvního symbolu.

• ve fázi H , přerozdělím řetěz do bucketů podle prvního 2^H symbolů.

→ na začátku fáze se řetěz v každém bucketu shodují v prvních H znacích na konci v prvních 2^H symbolech (buckety se podrozdělí)

→ fáze trvá $O(m)$ čas.

↓

využívám, že již znám rozdělení řetězů podle H -symbolů

Pro porovnání $T[i, m]$ a $T[j, m]$

se srovnají první H symbolů kde

použít informaci o bucketech

$T[i+H, m]$ a $T[j+H, m]$.

• skvělá implementace bere buckety

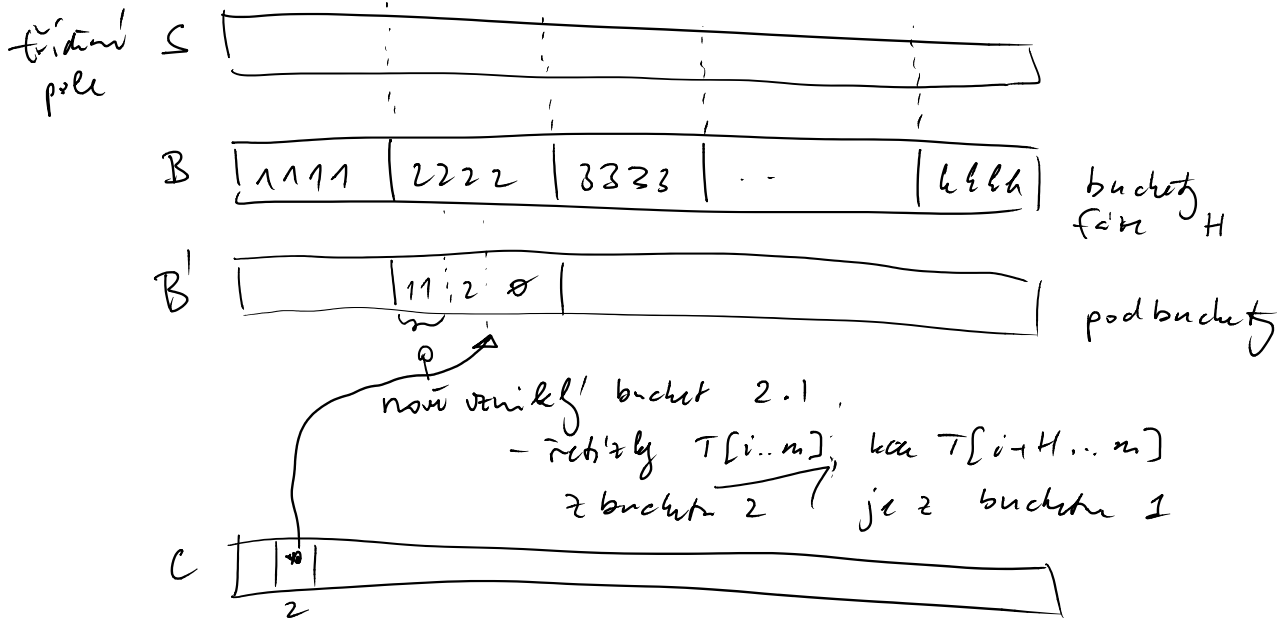
v lex-řádku, v rámci bucketu

jeho řetězku $T[i..m]$ po řetězku

a následně řetězku $T[i+H..m]$

(když $i-H > 0$)

v rámci jeho bucketu na začátek
do nově odděleného bucketu
[vše je uloženo v poli.]



nově vzniklý bucket 2.1.
 - řetězce $T[i..m]$, kde $T[i+H..m]$
 z bucketu 2 je z bucketu 1

ukazatel na ještě nepracovanou část
 v daném bucketu, kde vznikl
 nový podbucket

- na konci fáze přechází vzniklé buckety od 1, 2, ...

Sufixová stromy: trie sestavená z $T[1..n], T[2..n], \dots, T[m..n]$

- konceptuálně přimotává se ke sufixové pole,
 ale sufixové pole je kompaktnější, a hledání
 času obě řešení jsou srovnatelná.

- mnoho různých algoritmů na konstrukci suf. pole & stromů.
- komprimované trie, kde se místo řetězců
 na hranách používá odkazy na podřetězce
 v T potřebují pouze $O(m)$ místa.

Užití:

- 1) Hledání P v T
 - 2) zobrazení: Hledání P v T_1, T_2, \dots, T_k
 - 3) největší společný podřetězec (podinterval)

$$T_1 \text{ a } T_2 \dots O(|T_1| + |T_2|)$$
 (ignorujeme "log")
- a mnoho dalších

Samopravý: a seznamy

- množina prvků x_1, \dots, x_n
- chci reprezentovat spojité seznamem
- prvky p_1, p_2, \dots, p_n , kde p_i je pravděpodobnost, že Find bude hledat x_i .

otázka: Optimální uspořádání prvků v seznamu?
 → podle klesající pravděpodobnosti p_i

Bhava: $p_1 \geq p_2 \geq p_3 \dots \geq p_n$

očekávaný čas operace Find (x_i)

$$E[T] = \sum_{i=1}^n i \cdot p_i$$

s prvky p_i hledám x_i ~ musím
 tedy přečíst i prvků.

Problém: většinou nemáme p_1, p_2, \dots, p_n dopředu