


# Vyhledávání v textu:

$\Sigma :=$  abeceda  $\rightarrow$  množina všech znaků, chceme její konečnost i konstantnost  
základní i rozumně nulou, abychom mohli indexovat

$\Sigma^* :=$  řetězce  $\rightarrow$  množina všech řetězců

- |                                       |  |
|---------------------------------------|--|
| $\alpha, \beta, \gamma \dots$ řetězce | $\alpha \beta \dots$ zřetězení                           |
| $a, b, c \dots$ znaky                 | $\alpha[i]$ ... $i$ -tý znak $\alpha$ uvnitř             |
| $ \alpha $ ... délka řetězce          | $\alpha[i:j]$ ... podřetězec ( $i$ tam je, $j$ tam není) |
| $\epsilon$ ... prázdný řetězec        | $\alpha[:j]$ ... prefix                                  |
|                                       | $\alpha[i:]$ ... suffix                                  |

 Každý podřetězec je prefixem suffixu.

Problém: dostaneme slovo  $\sigma$ ,  $S := |\sigma| \rightarrow$  což je vstupní řetězec  
jehla  $\omega$ ,  $J := |\omega| \rightarrow$  což je hledaný řetězec  
výstup:  $\{i \mid \sigma[i:i+J] = \omega\}$

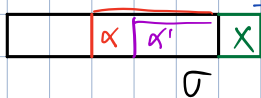
Triviální přístup:

$\Theta(S \cdot J) \rightarrow$  zkusím každý začátek a porovnám s jehlou následující

**Neefektivní:**  $\Theta(S) \cdot J$ , vstup:  $UOOUOOS$ , hledám  $UOOUO$   
Postupně probíráme, když nemáme shodu z hledáním, prostě zanedbáme.  
Zkusíme mi jen  $UOS$ , takže nic nemáme, přestože text  $UOOUOOS$  obsahuje.

## Implementační algoritmus:

Postupně čtem vstup a podle toho dopočítáváme výstup. Pokud  $\alpha$  není prefixem jehly, zanedbáme  $\alpha$  a



stav  $\alpha :=$  největší prefix jehly, který je suffixem slova

$\hookrightarrow$  jehly při čtení  $UOOUO|UO|U$   
 $\rightarrow$  zde  $UOOUO$   
 $\rightarrow$  zde  $UOOUO$   
 $\rightarrow$  zde bude  $UOUU$   
 $\rightarrow$  zde bude  $UOUU$  ( $UOOUU$  už není prefix jehly)

$\alpha'$  tak dlouho, než  $\alpha'$  je prefixem jehly

nový stav  $\rightarrow \epsilon$  (prázdný řetězec)

$\alpha'x$  suffix  $\sigma x$   
 $\hookrightarrow \alpha'$  suffix  $\sigma$   
 $\hookrightarrow \alpha'x$  prefix jehly  $\hookrightarrow \alpha'$  prefix jehly

$\left. \begin{array}{l} \epsilon \text{ to ho plyne: } |\alpha'| \leq |\alpha| \\ \alpha' \text{ je suffixem } \alpha \end{array} \right\}$

Předpokládáme posunuti:

$\rightarrow \alpha' \neq \alpha$

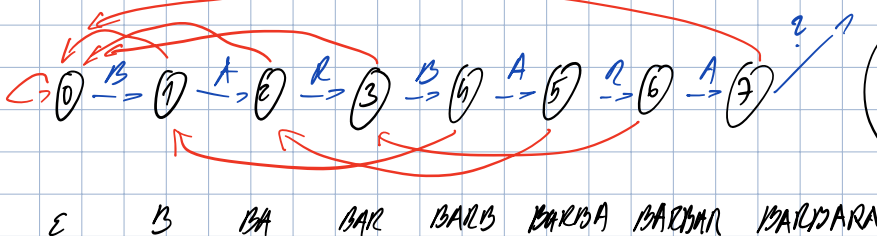
Pro slovo  $\alpha := z(\alpha) :=$  nejdelší  $\alpha'$  vlastní suffix  $\alpha$ , který je prefixem jehly

KMP algoritmus (Knuth, Morris, Pratt):

lze pospat pomocí vyhledávacího automatu.

Řetězec: BARBARA

stav:



pojištění proti overflow zpět -> nejbližší místo, kam se můžeme vrátit, abychom měli funkční prefix  
dopřední -> posunuté jehly

pole  $z[0..j]$

KMP:

Ukrok  $(i, x)$ :

Pokud  $jehla[i] \neq x$ :

Jeli  $i=0$ : vrátíme 0

$i = z[i]$

Vrátíme  $i+1$

Hledání jehly:

rychlost?:

Hleděj  $(\sigma)$

1.  $i=0$

2. Pro  $j=0 \dots |\sigma|-1$ :

3.  $i = \text{ukrok}(i, \sigma[j])$

4. Pokud  $i = |jehla|$ :

5. Ohlásím výsledek na  $j - |jehla| + 1$ .

Lemma: Hleděj  $(\sigma)$  má složitost  $\Theta(|\sigma n|)$ .

Důk: Hleděj je zjevně lineární, pouze uvnitř  $\text{ukrok}(i, x)$  může být problém.

# dopředných hran je  $\leq |\sigma n|$

# zpětných  $\leq$  # dopředných  $\leq |\sigma n|$

Celkem tedy nejvíce  $2|\sigma|$  hran projde.

Celý algoritmus je tedy  $\Theta(|\sigma n|)$

nejdelší vlastní suffix

☞ Pokud  $\alpha$  je slovo a spustíme automat na vstup  $\alpha[1..j]$ , vrátíme ve stavu  $z(\alpha)$

- pokud bychom našli jiný suffix než celou  $\alpha$

- zpětní hranu se dají konstruovat až za chvilu -> všechny předpoklady už budou vyhovovat.

Spustím na:

A  
AR  
ARB  
ARBA  
ARBAR  
ARBARA

stává nabit jen hole

ARBARA

# Konstrukce automatu:

1)  $z[0] = 0, z[1] = 0$

2)  $i = 0$

3) Pro  $j = 2, \dots, J:$

4)  $i = \text{Krok}(i, \text{jeht}[j-1])$

5)  $z[j] = i$

$\theta(j) + \theta(s)$  složitost hledj (s)

V: KMP najde všechny výskyty v čase  $\theta(J+S)$ .

Důležité informace.

Obecněji:

Jehtly:  $u_1 \dots u_n$

délka:  $J_1 \dots J_n$

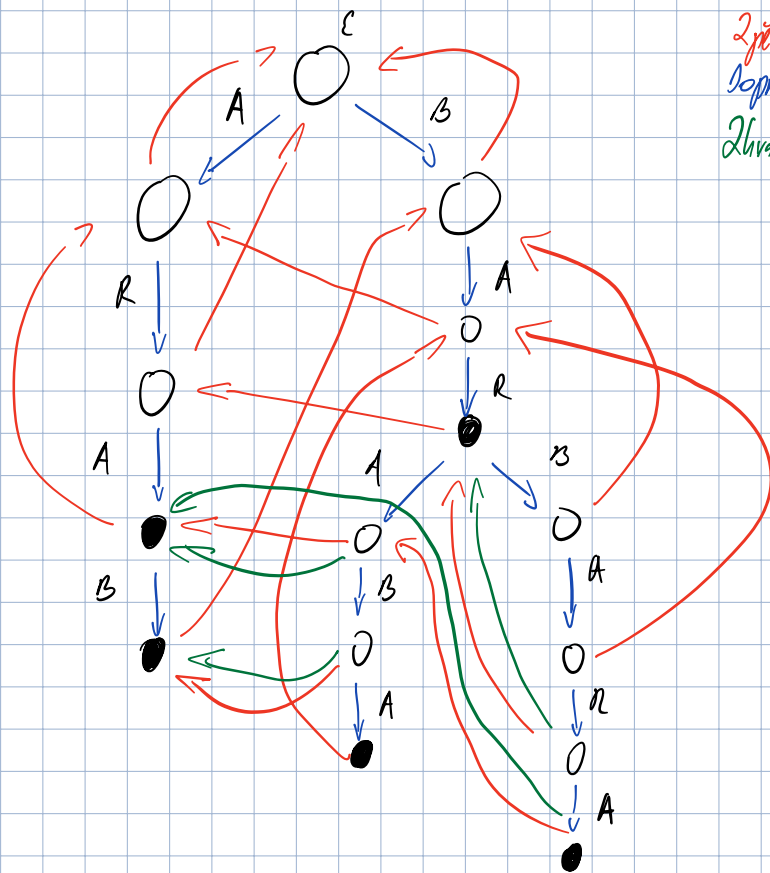
Seho  $\sigma$  délky  $S$

Chceme:  $\{(i,j) \mid u_i = \sigma[j:j+|u_i|]\}$

čas  $\theta(S + \sum J_i + V) \rightarrow \# \text{ výskytů}$

Jehtly:

- ARA
- BAR
- ARAB
- BARBARA
- BARABA



zpětní  
dopřední  
zlvatka  
do nejbližšího  
koncového bodu  
po zpětných  
hraničích