

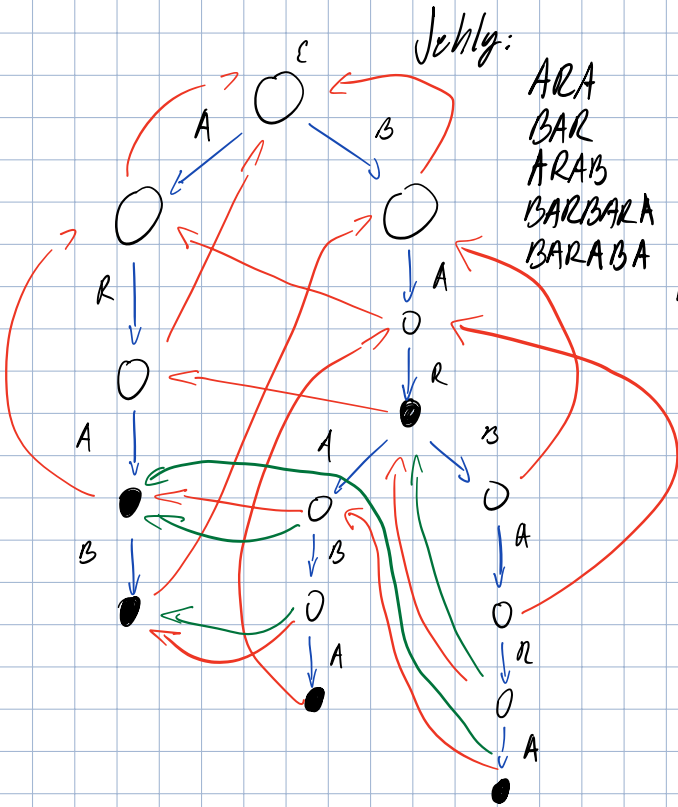
# KMP - obecněji (Aho-Corasick)

$i_1 \dots i_n :=$  jehly

$\sigma :=$  slovo

stav  $y :=$  prefixy všech jehel

hrany:   
 - dopředná  $\rightarrow$  prodlážením prefixu o jednu značku   
 - zpětná  $\rightarrow$  zkrácením prefixu na nejdelší vlastní suffix   
 - zkratková  $\rightarrow$  cesta po zpětných, jde na největší koncový stav - tedy literou jehly již obsahuje



## Reprezentace automatu:

stav  $o$  řešíme,  $o =$  kořen ( $c$ )

slovo  $(i) :=$  která jehla končí ve stavu  $i$

zpět  $(i) :=$  kam vede zpětná hrana

zkratka  $(i) :=$  kam vede zkratková hrana

Dopředn  $(i, x) :=$  kam vede dopředná hrana pro písmeno  $x$

*Vždy maximálně jedna od každé na jednom místě.*

## Krok $(s, x)$ :

*To je ze pro zadání  $x$   $o$   $x$  už není součástí řetězce*

Dokud **Dopředn  $(s, x) = \emptyset$** :  $\rightarrow$  pokud už není kam dál jít, oba sláhnou zpět

Dokud  $s =$  kořen: vrátíme  $s$ .  $\rightarrow$  pokud jsem již v kořenu

$s =$  zpět  $(s)$   $\rightarrow$  jímž jedu po zpětných

Vrátíme Dopředn  $(s, x)$   $\rightarrow$  už existuje příslušná dopředná hrana

## Hledj $(\sigma)$ :

$s =$  kořen

Pro  $i = 0 \dots |\sigma| - 1$ :

$s =$  Krok  $(s, \sigma[i])$

$t = s$  - temp. stav

Dokud  $t \neq \emptyset$ :

Jeli slovo  $(t) \neq \emptyset$ :

Hlášíme výslyt...

$t =$  zkratka  $(t)$   $\rightarrow$  takže funkce jen dokud je def.

*Lemma: Hledj bž. v čase  $O(|\sigma| + \# \text{ výslytů})$*

$\#$  dopředných je nejvýše tolik co  $|\sigma|$

$\#$  zpětných je nejvýše tolik co dopředných

*tedy  $O(|\sigma|)$*

$\#$  výslytů je tolik, kolikrát se ohlásí.

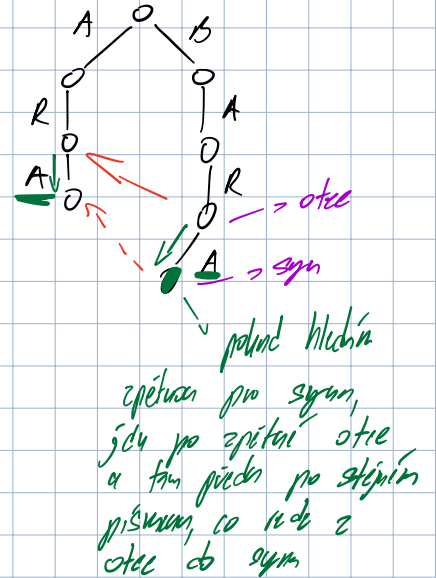
$s$  tím, že poprvé se může projít slyšelnou bez ohlášení...

## Konstrukce automatu:

Myslelím: Stejně jako u UMP, myslí každá hrana počítat všechny jehly paralelně po jednom písmenu a strom bude tvořit po restách.

## Konstrukce ( $i_0 - i_n$ ):

- 1) Vybereme tři pro  $i_1 - i_2 \rightarrow 2$  toho máme dopředu
- 2) Zpětní (kořen) =  $\emptyset$ , Zhratka (kořen) =  $\emptyset$
- 3) Zpětní ( $\forall$  sym kořene) = kořen, Zhratka ( $\forall$  sym kořene) =  $\emptyset$
- 4)  $F$  = frata se sym kořene:
- 5) Dokud  $F \neq \emptyset$ :
- 6)  $V$  = deque(F)
- 7) Pro  $s$  sym  $v$ :
- 8)  $zpět(s) = \text{Urok}(zpět(v), \text{písmeno na hraně } vs)$
- 9)  $F.Add(s)$
- 10) Pokud Slovo( $zpět(s)$ )  $\neq \emptyset$ :
- 11)  $Zhratka(s) = zpět(s)$
- 12) Jinak  $Zhratka(s) = Zhratka(zpět(s))$ .



$\rightarrow$  Využití funkce, že se chodí po zpětných a předstí je málo nalezené a pospojované

## Lemna: Konstrukce běží v čase $O(\epsilon \cdot i)$

Bca Urok je vyhledávání tří asympt. lineární. Zbytek je BFS.

Jehly procházíme paralelně a procházíme jednou jehlou je lineární.

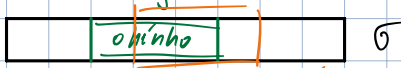
A tedy i celý paralelní přístup je lineární.  $\rightarrow$  Dokonce často stav počítán jen jednou a u dalších jehel ho nemusíme počítat.

Věta: Nalezení všech výskytů jehel v seně trvá  $O(\epsilon \cdot i + |S| + \# \text{ výskytů})$ .

Důkazy jsou spojené lemmat.

# Robinův - Karpiův alg.

$$h(x_0 - x_{j-1}) \rightarrow \partial \in \mathbb{Z}_H$$



musím umět přepočítat  
hrušky v konstantním  
časě po posunutí okénka  
jinak bych nic neudělal!

$$h(x_0 - x_{j-1}) := (x_0 r^{j-1} + x_1 r^{j-2} + \dots + x_{j-1} r^0) \bmod H$$

posunutí okénka

$$h(x_1 - x_j) =$$

$$(x_1 r^{j-1} + x_2 r^{j-2} + \dots + x_j r^0) \bmod H$$

Dobrá?

$x_0$  zmizelo,  $x_j$  přibýlo  
ostatní se jen vynásobilo  $r$ -kem

Řešení:

$$h(x_1 - x_j) = h(x_0 - x_{j-1}) \cdot r - x_0 r^j + x_j$$

A tabule je  
konstantní  
čas  $O(1)$

## RL alg.

- 0) zvolíme  $r$  z tělesa náhodně
- 1)  $c = h(\text{jablko})$ ,  $a = h(\sigma[1:j])$ , spočítám  $r^j$
- 2) pro  $i = 0, \dots, |\sigma| - j$ :
- 3) Pokud  $a = c$  and  $\sigma[i:i+j] = c$
- 4) Hlásím výslyt jablka  $i$  na pozici  $i$ .
- 5)  $a = (a \cdot r - \sigma[i] \cdot r^j + \sigma[i+j]) \bmod H$ .

### Složitost:

- řešení + počítání hrušek  $= O(|\sigma|)$

- skutečné výslyty  $= O(j \cdot v)$  → výslyty

- falešné výslyty =

↳ pro ideální (obdobně náhodnou)  $h$ ,

$$P(\text{falešného výslytu}) = \frac{1}{H}$$

$$\Rightarrow \text{průměrný čas} = O\left(\frac{|\sigma| \cdot j}{H}\right)$$

$\frac{j}{H}$  pro jednu  
dílnku.  
Celkem  $|\sigma|$   
dílnek.

Pokud  $H \geq j$ :  $O(|\sigma|)$

Pozor: skutečný výslyt  
je hodně drahý.  
Vhodné pro normalizované jablko.

Mějme  $P(x) := p_0 x^0 + p_1 x^1 + \dots + p_{n-1} x^{n-1}$  nad tělesem

Lemna: Pokud  $x_1, \dots, x_n$  jsou všechny kořeny polynomu  $P$ , pak  $P(x) = (x-x_1) \cdot (x-x_2) \cdot \dots \cdot (x-x_n) \cdot Q(x)$

kde  $Q(x)$  je polynom bez kořenů.

Důsledek: Polynom stupně  $d$  má nejvýše  $d$  kořenů

Uděl by  $k$  bylo větší než stupeň  $P$ , vyšel by spor.

Lemma: Necht'  $P$  a  $Q$  jsou polynomy stupně  $< n$ ,  $x_1, \dots, x_n$  náravným různými čísly t.č.  $\forall i: P(x_i) = Q(x_i)$ .  
Pak  $P \equiv Q$ .

$R := P - Q$ , stupeň opět  $< n$  t.č.  $\forall i: R(x_i) = 0$   $\rightarrow$  Jediný takový polynom je nulový polynom.

$$R \equiv 0 \Rightarrow P \equiv Q$$

Tedy nyní falešné výsledky doplníme:

$$P(r) = h(i)$$

$$Q(r) = h(\sigma[i:i+j])$$

- Ptáme se, že pokud jsou polynomy různé, tak pro konkrétní  $r$  jistě je pravděpodobnost, že se budou rovnat.

Tedy  $P[P(r) \equiv Q(r)] \leq \frac{j}{H} \rightarrow \frac{1}{H}$  opět z velikosti množ. prvků,  $\leq j$ -líst, protože to je také pozice, kde se dva polynomy mohou rovnat, aby stále byly identické.

Tedy  $H$  by muselo být alespoň  $j^2$  velká.

Celková časová složitost je tedy  $O\left(\frac{\kappa \cdot j^2}{H}\right) \rightarrow$  tedy pro  $H \geq j^2$  je to  $O(\kappa)$

$O(\kappa \cdot j) :=$  ověření jehly v seně

$\cdot O\left(\frac{j}{H}\right)$  počet falešných výsledků (kromě toho správného).