

## Úplně indexování textu.

### Lyndon decomposition

- primitivní slovo & ostatě nejmenší lexicografické rotace

cbcbbaabacacacaca - |c|bc|b|b|aabaac|aaac|a

LZ - decomposition - používáno v LZ kompresi

Border Array - obsahuje délku největšího borderu na každé pozici

- pohled ze dalšího hlediska, takže na  $i = \beta(\text{current})$

→ pohled menšího pokračování, začínají 0

$$[\beta[1] = 0]$$

Suffix array - nové symboly  $\$ < \Sigma < \#$

string  $T = \# \text{abaabab} \$$

i	-1	0	1	2	3	4	5	6	7
$T[i]$	#	a	b	a	a	b	a	b	\$
$SA[i]$	7	2	5	0	3	6	1	4	-1
$LCP[i]$	0	0	1	2	3	0	1	2	0

starting pozice lexicograficky seřazených suffixů

největší společné ~~sub~~ prefixy

### LCP Search

$x = \text{aba}$  - hledaný string

algorithmus:  $d = -1$   
 $f = 7$

$i = \lfloor (d+f)/2 \rfloor = 3$  - průměrná pole

$L3 = T[SA[3] \dots 13] = T[3 \dots 6] = \text{abab}$

$l = LCP(L3, x) = 3 \rightarrow \text{máme} - \text{našli jsme}$

- přes LCP kolem pozice 3 máme největší ostatní výskyt

- musí být  $\geq \text{len}(x)$

- pokud bych neměli, porovnáváme  $L3[1] \text{ s } x[1]$

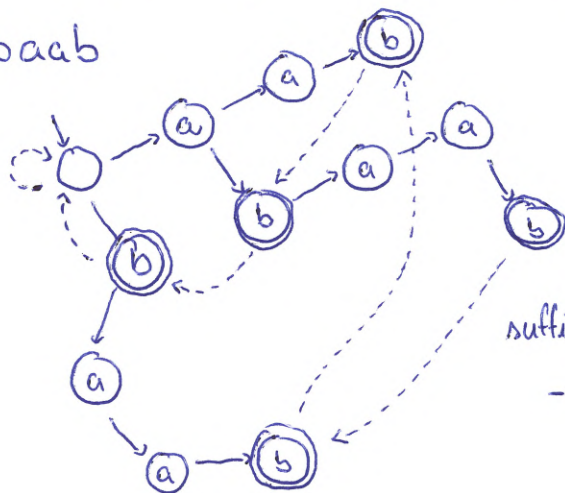
→ pokud je  $x[1]$  menší, pokračují v menší polce (tj  $f = i$ )

→ v opačném případě se vyšší ( $d = i$ )

Suffix Trie - ~~skem~~ skem všech možných suffixů

- stav pro každý možný mezislov, label hranu vždy jeden symbol

$T = abaaab$



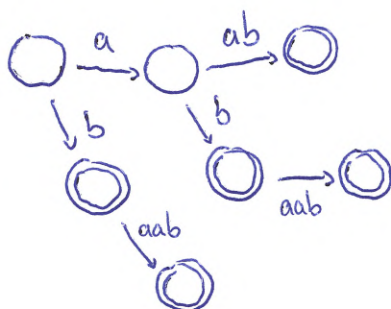
- uzel s větvením: fork

- koncový uzel vždy má konci suffixu

suffixové linky

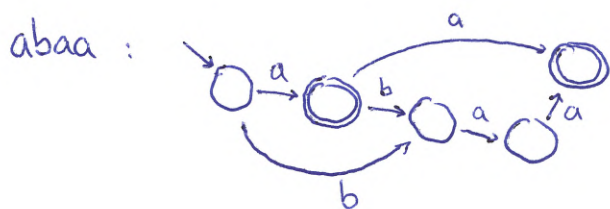
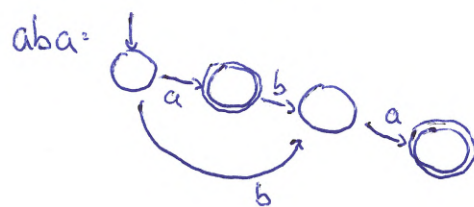
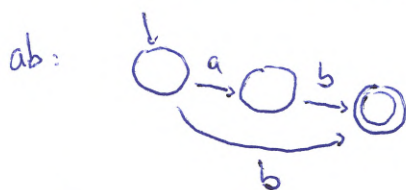
- měříteji do uzel se suffixem o 1 znakem

Suffix Tree - trie with nonfinal nodes of degree 1 deleted

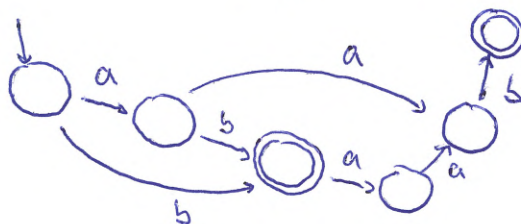


Suffix automaton

- minimální automat přijímající všechny suffixy



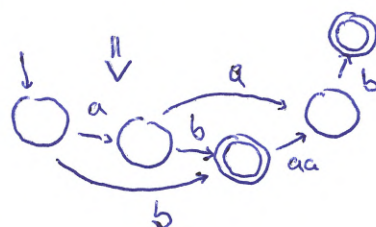
abaab:



- 15 každé iteraci "vytvoří koncový stav" ze předešlého stavu

- pokud se sejde se každé předešlé hranou koncový stav → de

- 15 opakování případě je nutný fork



Compact suffix automaton  
[implementation]