

## Suffix Automata

MAXIME CROCHEMORE

King's College London

Maxime.Crochemore@kcl.ac.uk  
<http://www.dcs.kcl.ac.uk/staff/mac/>

## Implementation of indexes



Implementation with efficient data structures

- ★ **Suffix Trees**  
digital trees, PATRICIA tree (compact trees)
- ★ **Suffix Automata or DAWG's**  
minimal automata, compact automata

Implementation with efficient algorithm

- ★ **Suffix Arrays**  
binary search in the ordered list of suffixes

## Suffixes

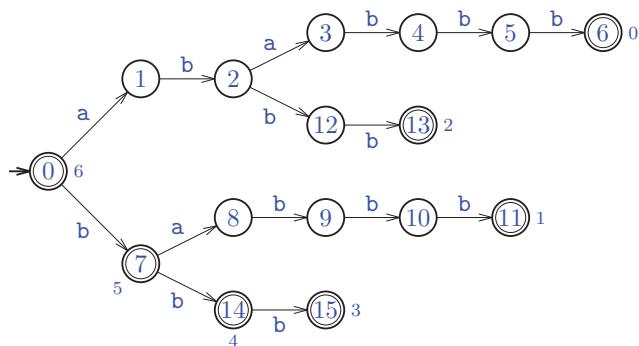
Text  $y \in \Sigma^*$

- ★  $\text{Suff}(y) = \text{set of suffixes of } y,$
- ★  $\text{card } \text{Suff}(y) = |y| + 1$
- ★  $\text{Suff}(\text{ababbb})$

$i$	0	1	2	3	4	5	
$y[i]$	a	b	a	b	b	b	
							position
	a	b	a	b	b	b	0
	b	a	b	b	b		1
	a	b	b	b			2
	b	b	b				3
	b	b					4
	b						5
	$\varepsilon$						(empty string)

## Trie of suffixes

- ★  $\mathcal{T}(y) = \text{digital tree which branches are labeled by suffixes of } y$   
= tree-like deterministic automaton accepting  $\text{Suff}(y)$
- ★ **Nodes**  
identified with factors (subwords) of  $y$
- ★ **Terminal nodes**  
identified with suffixes of  $y$ , output = position of the suffix
- ★ **Suffix trie of ababbb**

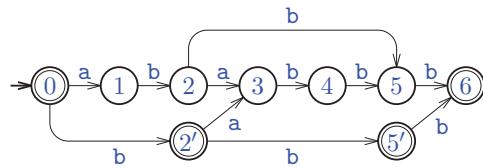


## Suffix Automaton

Text  $y \in \Sigma^*$  of length  $n$

$\mathcal{A}(y)$  = minimal deterministic automaton accepting  $Suff(y)$

- ★ **Minimization** of the trie of suffixes



- ★ States are classes of factors (subwords) of  $y$

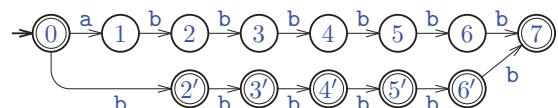
- ★ **Size:**

$$n + 1 \leq \#states \leq 2n - 1$$

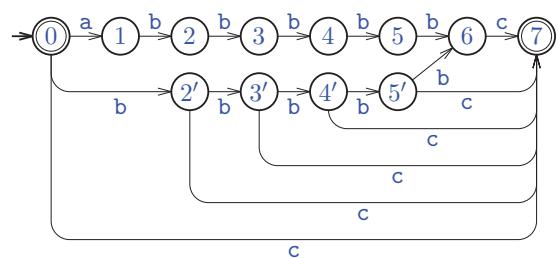
$$n \leq \#arcs \leq 3n - 4$$

## Maximal size

- ★ **Maximal number of states**



- ★ **Maximal number of arcs**

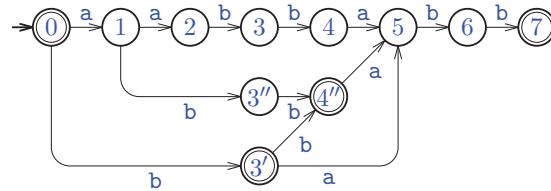


## Suffix link

★ **Function  $f_y$ , suffix link**

let  $p = \text{TARGET}(\text{initial}[\mathcal{A}], v)$ ,  $v \in \Sigma^+$

$f_y(p) = \text{TARGET}(\text{initial}[\mathcal{A}], u)$ , where  $u$  is the longest suffix of  $v$  occurring in a different right context



- ★  $f[1] = 0, f[2] = 1, f[3] = 3'', f[3''] = 3', f[3'] = 0,$   
 $f[4] = 4'', f[4''] = 3', f[5] = 1, f[6] = 3'', f[7] = 4''.$

★ **Suffix path**

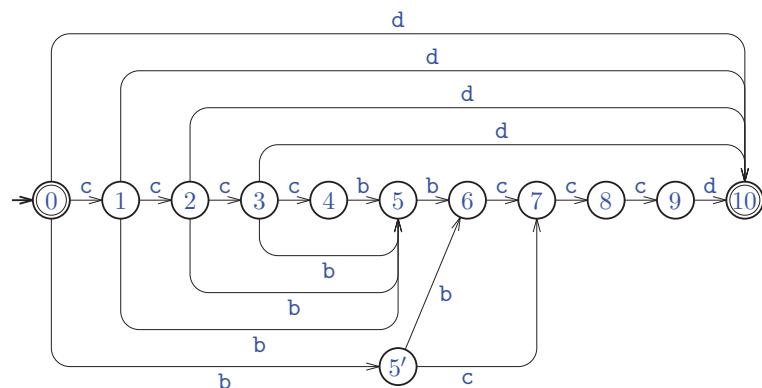
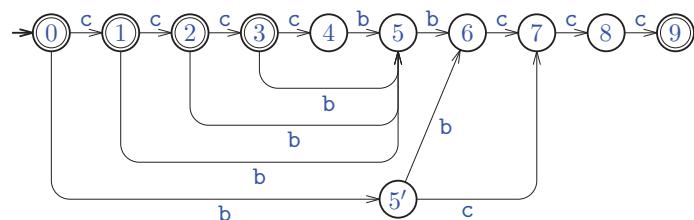
example for state 7:  $\langle 7, 4'', 3', 0 \rangle$ , sequence of terminal states

★ **Use**

same but more efficient than suffix link in suffix trees

## Construction—one step (1)

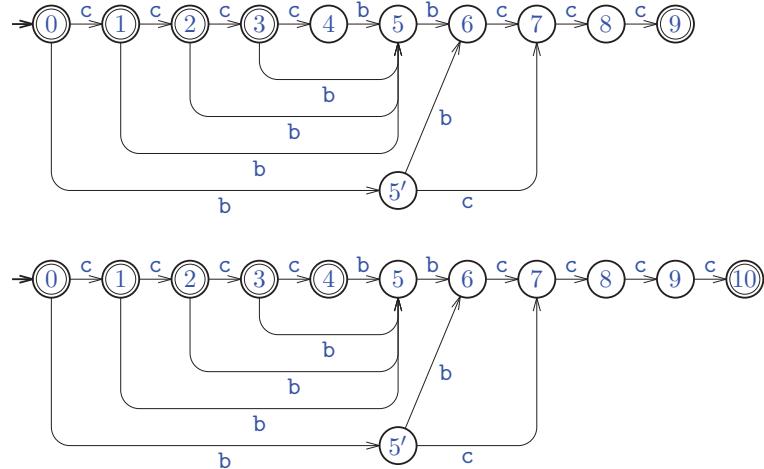
★ **From  $\mathcal{A}(ccccbbccc)$  to  $\mathcal{A}(ccccbbcccd)$**



- ★ New arcs from states of the suffix path  $\langle 9, 3, 2, 1, 0 \rangle$ .

## Construction—one step (2)

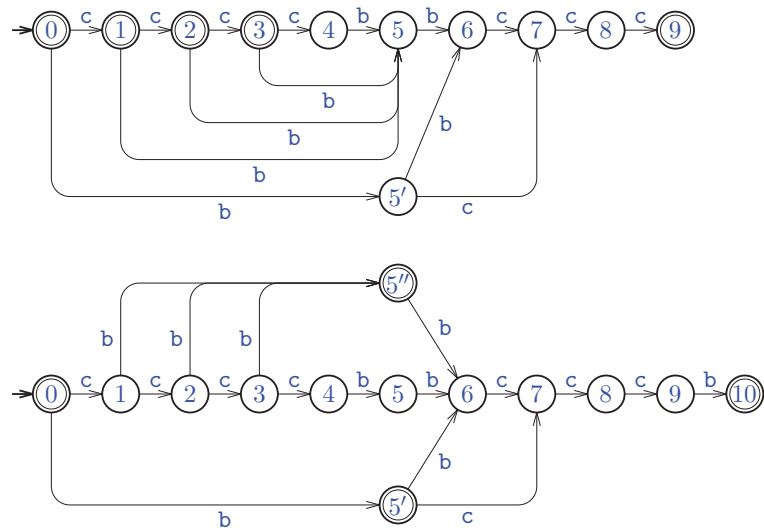
- ★ From  $\mathcal{A}(\text{ccccbbccc})$  to  $\mathcal{A}(\text{ccccbbccccc})$



- ★ Link  $3 = f[9]$  and solid arc  $(3, c, 4)$  (not a shortcut)  
then,  $f[10] = \text{TARGET}(3, c) = 4$  that becomes a terminal state

## Construction—one step (3)

- ★ From  $\mathcal{A}(\text{ccccbbccc})$  to  $\mathcal{A}(\text{ccccbbcccb})$



- ★ Link  $3 = f[9]$ , non-solid arc  $(3, b, 5)$ , **cccb** suffix but **ccccb** not  
state 5 is cloned into  $5'' = f[10] = f[5], f[5''] = 5'$   
arcs  $(3, b, 5), (2, b, 5)$  et  $(1, b, 5)$  are redirected onto  $5''$

## Construction—main

**SUFFIX-AUTO**( $y, n$ )

```

1   $M \leftarrow \text{NEW-AUTOMATON}()$ 
2   $L[\text{initial}[M]] \leftarrow 0$ 
3   $F[\text{initial}[M]] \leftarrow \text{NIL}$ 
4   $\text{last}[M] \leftarrow \text{initial}[M]$ 
5  for each letter  $a$  of  $y$  sequentially do
6       $\triangleright$  Extension of  $M$  by the letter  $a$ 
7       $\text{EXTENSION}(a)$ 
8   $p \leftarrow \text{last}[M]$ 
9  do  $\text{terminal}[p] \leftarrow \text{TRUE}$ 
10      $p \leftarrow F[p]$ 
11 while  $p \neq \text{NIL}$ 
12 return  $M$ 
```

## Construction—extension

**EXTENSION**( $a$ )

```

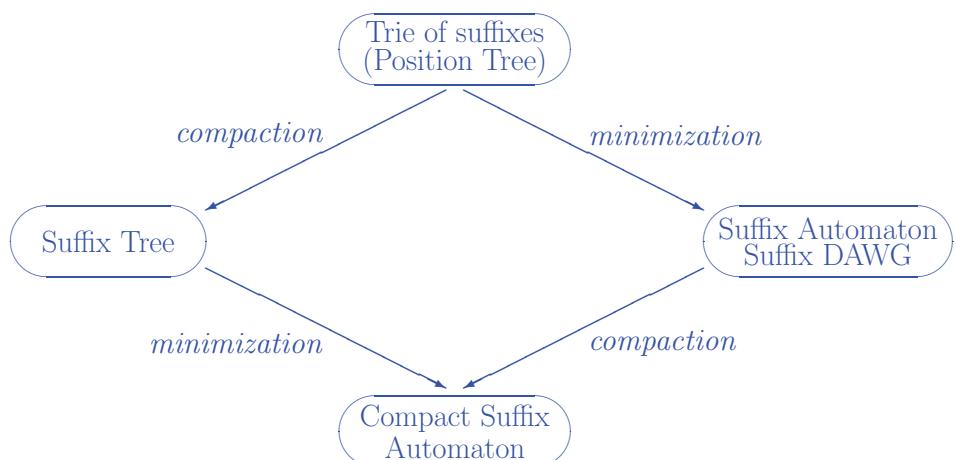
1   $\text{new\_last} \leftarrow \text{NEW-STATE}()$ 
2   $L[\text{new\_last}] \leftarrow L[\text{last}[M]] + 1$ 
3   $p \leftarrow \text{last}[M]$ 
4  do  $\text{Adj}[p] \leftarrow \text{Adj}[p] \cup \{(a, \text{new\_last})\}$ 
5       $p \leftarrow F[p]$ 
6  while  $p \neq \text{NIL}$  and  $\text{TARGET}(p, a) = \text{NIL}$ 
7  if  $p = \text{NIL}$  then
8       $F[\text{new\_last}] \leftarrow \text{initial}[M]$ 
9  else  $q \leftarrow \text{TARGET}(p, a)$ 
10     if  $(p, a, q)$  is solid, i.e.  $L[p] + 1 = L[q]$  then
11          $F[\text{new\_last}] \leftarrow q$ 
12     else  $\text{clone} \leftarrow \text{NEW-STATE}()$ 
13          $L[\text{clone}] \leftarrow L[p] + 1$ 
14         for each pair  $(b, q') \in \text{Succ}[q]$  do
15              $\text{Adj}[\text{clone}] \leftarrow \text{Adj}[\text{clone}] \cup \{(b, q')\}$ 
16              $F[\text{new\_last}] \leftarrow \text{clone}$ 
17              $F[\text{clone}] \leftarrow F[q]$ 
18              $F[q] \leftarrow \text{clone}$ 
19             do  $\text{Adj}[p] \leftarrow \text{Adj}[p] \setminus \{(a, q)\}$ 
20                  $\text{Adj}[p] \leftarrow \text{Adj}[p] \cup \{(a, \text{clone})\}$ 
21                  $p \leftarrow F[p]$ 
22             while  $p \neq \text{NIL}$  and  $\text{TARGET}(p, a) = q$ 
23      $\text{last}[M] \leftarrow \text{new\_last}$ 
```

## Operations on indexes

Text  $y$  of length  $n$

- ★ **Index implemented by suffix tree or suffix automaton of  $y$**   
memory space  $O(n)$ , construction time  $O(n \times \log \text{card } \Sigma)$
- ★ **String matching**  
searching  $y$  for  $x$  of length  $m$ : time  $O(m \times \log \text{card } \Sigma)$   
number of occurrences of  $x$  in  $y$ : same complexity after  $O(n)$  preprocessing
- ★ **All occurrences**  
finding all occurrences of  $x$  in  $y$ : time  $O(m \times \log \text{card } \Sigma) + |\text{output}|$
- ★ **Repetitions**  
computing a longest factor of  $y$  occurring at least  $k$  times: time  $O(n)$
- ★ **Marker**  
computing a shortest factor of  $y$  occurring exactly once: time  $O(n)$

## Saving space

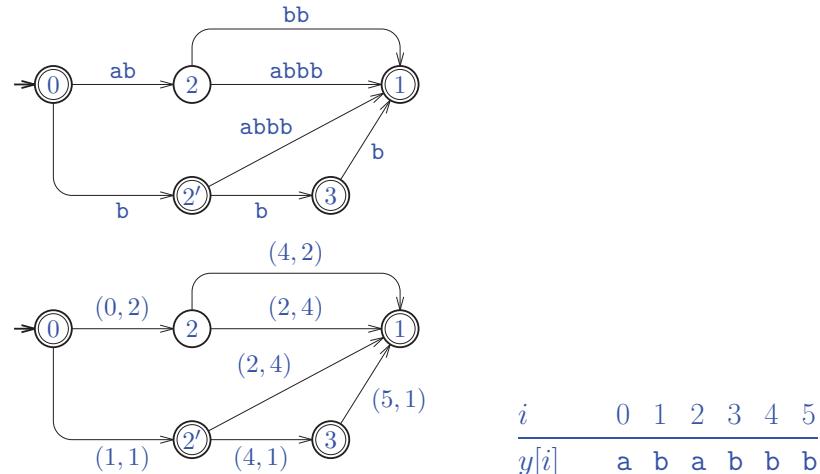


## Compact Suffix Automaton

Text  $y \in \Sigma^*$  of length  $n$

$\mathcal{A}^c(y)$  = compact minimal automaton accepting  $Suff(y)$

★ **Compaction** of  $\mathcal{A}(y)$ , or **minimization** of  $\mathcal{S}(y)$



★ **Linear size:**  $O(n)$

## Direct construction of CSA

- ★ Similar to both
  - Suffix Tree construction
  - Suffix Automaton construction
- ★ Sequential addition of suffixes in the structure from the longest to the shortest
- ★ Used features:
  - “slow-find” and “fast-find” procedures
  - suffix links
  - solid and non-solid arcs
  - state splitting
  - re-directions of arcs
- ★ **Complexity:**  $O(n \log \text{card } \Sigma)$  time,  $O(n)$  space  
50% saved on space of Suffix Automaton