

## Dictionary Matching

MAXIME CROCHEMORE

King's College London

[Maxime.Crochemore@kcl.ac.uk](mailto:Maxime.Crochemore@kcl.ac.uk)

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

## Matching several strings

- ★ **Dictionary:** set of finite strings,  $X = \{x_0, x_1, \dots, x_{k-1}\}$   
(empty string not in  $X$ )
- ★ **Matching:** locate occurrences of the strings in any text  $y$

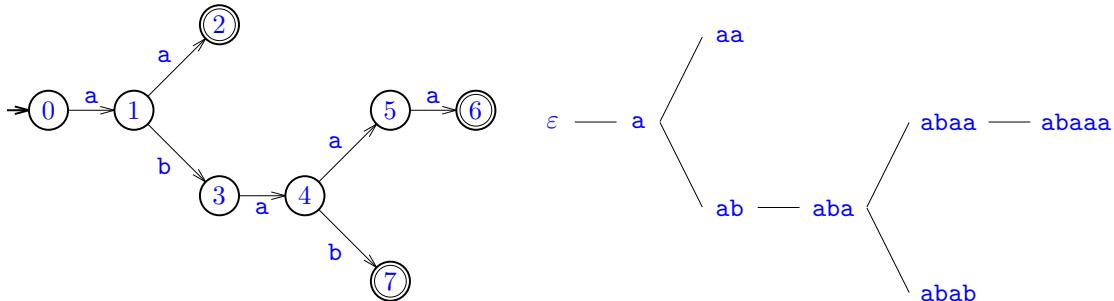
Note that each position on  $y$  can be the position of several strings, yielding a possible quadratic output (e.g.  $X = \{\text{a}, \text{aa}, \dots, \text{a}^k\}$  and  $y = \text{a}^n$ )

- ★ **Output:** list of positions on  $y$  that are end-positions of some string in  $X$
- ★ **Standard method:** based on the Dictionary Matching Automaton (Aho-Corasick automaton), an extension of the String Matching Automaton
- ★ **Other method:** extension of the right-to-left window scanning strategy (Boyer-Moore technique)

## Trie of the dictionary

- ★ **Trie:**  $\mathcal{T}(X)$ , digital tree whose branches are labelled by strings of  $X$   
As an automaton, the language it accepts is  $X$

- ★ **Example :**  $\mathcal{T}(\{aa, abaaa, abab\})$



- ★ Nodes are all prefixes of strings in  $X$
- ★  $\mathcal{T}(X)$  is the basis of the Dictionary Matching Automaton

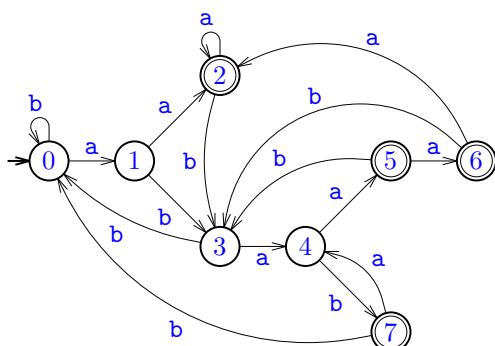
## Dictionary Matching Automaton

- ★ **Dictionary Matching Automaton,  $\mathcal{D}(X)$**

it accepts the language  $A^*X$  and is defined by:

- set of states is  $\text{Pref}(X)$ ; initial state is the empty string
- set of terminal states is  $\text{Pref}(X) \cap A^*X$
- arcs are of the form  $(u, a, h(ua))$   
where  $h(ua)$  = longest suffix of  $ua$  that belongs to  $\text{Pref}(X)$

- ★ **Example :**  $\mathcal{D}(\{aa, abaaa, abab\})$  with alphabet  $A = \{a, b\}$



## Searching with DMA

★ **Searching**

**DET-SEARCH-BY-FAILURE( $X, y$ )**

```

1   $M \leftarrow \text{DMA-BY-FAILURE}(X)$ 
2   $r \leftarrow \text{initial}[M]$ 
3  for each letter  $a$  of  $y$ , sequentially do
4     $r \leftarrow \text{TARGET}(, )(r, a)$ 
5     $\text{OUTPUT-IF}(\text{terminal}[r])$ 

```

★ **Implementation** of  $\mathcal{D}(\{\text{aa, abaaa, abab}\})$  with a transition table

Next state table	0	1	2	3	4	5	6	7
<b>a</b>	0	2	2	4	5	6	2	4
<b>b</b>	1	3	3	0	7	3	3	0

★ **Searching time:**  $O(|y|)$

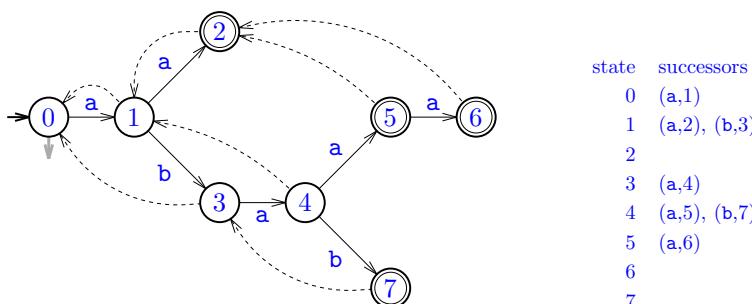
$j$	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	...	
$y[j]$	c	d	a	b	b	a	b	a	a	b	a	b	a	b	a	a	...		
state	0	0	0	1	3	0	1	3	4	<b>5</b>	3	4	<b>7</b>	4	<b>7</b>	0	1	<b>2</b>	...

## Implementation by failure function

- ★ **Aim:** reduction of the implementation to  $O(\Sigma\{|x| : x \in X\})$  independent of the alphabet

- ★ **Failure function** ( $u$  nonempty string)  
 $f(u) =$  the longest proper suffix of  $u$  that belongs to  $\text{Pref}(X)$

- ★ **Implementation** of  $\mathcal{D}(\{\text{aa, abaaa, abab}\})$  with successor lists and  $f$



★ **Searching**

$j$	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	...	
$y[j]$	c	d	a	b	b	a	b	a	a	b	a	b	a	b	a	a	...		
state	0	0	0	1	3	0,0	1	3	4	<b>5</b>	2,1,3	4	<b>7</b>	3,4	<b>7</b>	3,0,0	1	<b>2</b>	...

## Searching with failure function

- ★ **Searching**

**DET-SEARCH-BY-FAILURE( $X, y$ )**

```

1   $M \leftarrow \text{DMA-BY-FAILURE}(X)$ 
2   $r \leftarrow \text{initial}[M]$ 
3  for each letter  $a$  of  $y$ , sequentially do
4     $r \leftarrow \text{TARGET-BY-FAILURE}(r, a)$ 
5    OUTPUT-IF( $\text{terminal}[r]$ )

```

- ★ **Target:** next state function using successors lists

**TARGET-BY-FAILURE( $p, a$ )**

```

1  while  $p \neq \text{NIL}$  and  $\text{TARGET}(p, a) = \text{NIL}$  do
2     $p \leftarrow \text{fail}[p]$ 
3  if  $p = \text{NIL}$  then
4    return  $\text{initial}[M]$ 
5  else return  $\text{TARGET}(p, a)$ 

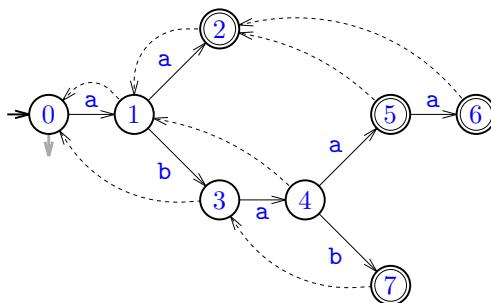
```

- ★ **Running time:**  $O(|y| \times \log \text{card } A)$

## Computing the failure links

- ★  $f(ua) = \text{TARGET-BY-FAILURE}(f(u), a)$

- ★ **Computation via a width-first traversal of the trie**



- ★  $\text{fail}[0], \text{fail}[1], \text{fail}[2], \text{fail}[3],$  and  $\text{fail}[4]$  already computed

- ★ **Computing  $\text{fail}[5]$  from 4:**  $\delta(4, a) = 5$ ,  $\text{fail}[4] = 1$ , and  $\delta(1, a) = 2$  gives  $\text{fail}[5] = 2$

**Note:** since state 2 is terminal, state 5 becomes terminal

- ★ **Computing  $\text{fail}[6]$  from 5:**  $\delta(5, a) = 6$ ,  $\text{fail}[5] = 2$ , and  $\delta(2, a)$  undefined but  $\text{TARGET-BY-FAILURE}(2, a) = 2$  gives  $\text{fail}[6] = 2$

## Computing the failure links—algorithm

- ★ Failure links  $\text{fail}[]$  implements the failure function  $f$

DMA-BY-FAILURE( $X$ )

```

1   $M \leftarrow \text{TRIE}(X)$ 
2   $\text{fail}[\text{initial}[M]] \leftarrow \text{NIL}$ 
3   $F \leftarrow \text{EMPTY-QUEUE}()$ 
4   $\text{ENQUEUE}(F, \text{initial}[M])$ 
5  while non QUEUE-IS-EMPTY( $F$ ) do
6     $t \leftarrow \text{DEQUEUED}(F)$ 
7    for each pair  $(a, p) \in \text{Succ}[t]$  do
8       $r \leftarrow \text{TARGET-BY-FAILURE}(\text{fail}[t], a)$ 
9       $\text{fail}[p] \leftarrow r$ 
10   if  $\text{terminal}[r]$  then
11      $\text{terminal}[p] \leftarrow \text{TRUE}$ 
12    $\text{ENQUEUE}(F, p)$ 
13 return  $M$ 

```

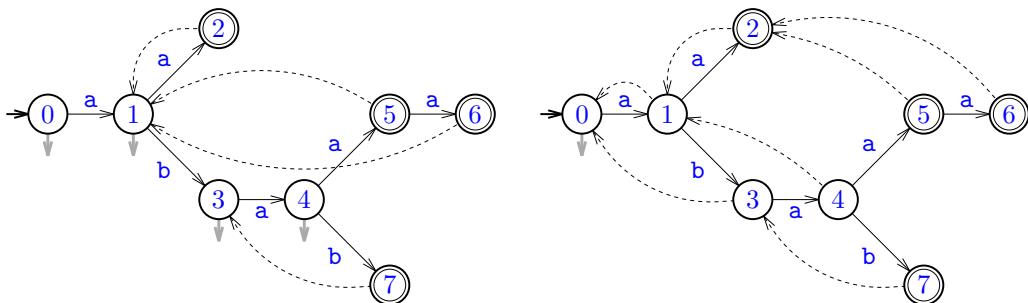
## Optimized failure links

- ★ Next set of state  $u \in \text{Pref}(X)$ :
- $$\text{Next}(u) = \{a : a \in A, ua \in \text{Pref}(X)\}$$

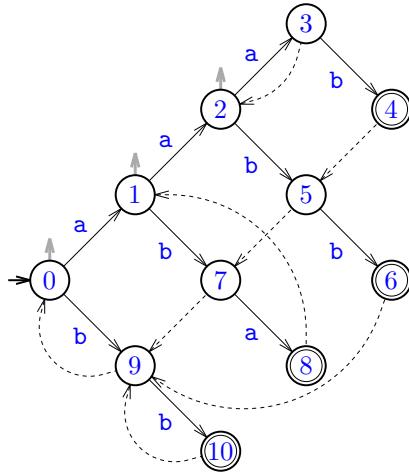
- ★ New link defined by function  $f'$ : for  $u$  nonempty  $f'(u) = f^k(u)$   
where  $k = \min\{\ell : \text{Next}(f^\ell(u)) \not\subseteq \text{Next}(u)\}$

- ★ Optimized link

Non optimized



- ★ **Delay:** maximum time spent on a letter of  $y$   
It is  $\max\{|x| : x \in X\}$  even with the optimized links
- ★ **Example:**  $X = \{aaab, aabb, aba, bb\}$



★  $L(m) = \{a^{m-1}b\} \cup \{a^{2k-1}ba : 1 \leq k < \lceil m/2 \rceil\} \cup \{a^{2k}bb : 0 \leq k < \lfloor m/2 \rfloor\}$