

## Searching a list of strings

MAXIME CROCHEMEORE

King's College London

Maxime.Crochemore@kcl.ac.uk

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

## Searching problem

### ★ Input

- a list  $L$  of  $n$  strings of  $\Sigma^*$  stored in increasing lexicographic order in a table:  $L_0 \leq L_1 \leq \dots \leq L_{n-1}$
- a string  $x \in \Sigma^*$  of length  $m$ .

### ★ Simple problem

find

- either  $i$ ,  $-1 < i < n$ , with  $x = L_i$  if  $x$  occurs in  $L$ ,
- or  $d$  and  $f$ ,  $-1 \leq d < f \leq n$ , that satisfy  $d + 1 = f$  and  $L_d < x < L_f$  otherwise.

### ★ Interval

find  $d$  and  $f$ ,  $-1 \leq d < f \leq n$ , with:  
 $d < i < f$  if and only if  $x$  prefix of  $L_i$ .

## Example

### ★ List $L$

$$L_0 = \text{a a a b a a}$$

$$L_1 = \text{a a a b b}$$

$$L_2 = \text{a a b b b b}$$

$$L_3 = \text{a b}$$

$$L_4 = \text{b a a a}$$

$$L_5 = \text{b b}$$

### ★ Search

$$x = \text{a a a b b} \longrightarrow 1$$

$$x = \text{a a b a} \longrightarrow (1, 2)$$

### ★ Interval

$$x = \text{a a} \longrightarrow (-1, 3)$$

## Searching algorithm

```
SIMPLE-SEARCH( $L, n, x, m$ )
1    $d \leftarrow -1$ 
2    $f \leftarrow n$ 
3   while  $d + 1 < f$  do           ▷ Invariant:  $L_d < x < L_f$ 
4        $i \leftarrow \lfloor (d + f)/2 \rfloor$ 
5        $\ell \leftarrow |\text{lcp}(x, L_i)|$ 
6       if  $\ell = m$  and  $\ell = |L_i|$  then
7           return  $i$ 
8       elseif  $(\ell = |L_i|)$  or  $(\ell \neq m \text{ and } L_i[\ell] < x[\ell])$  then
9            $d \leftarrow i$ 
10      else  $f \leftarrow i$ 
11      return  $(d, f)$ 
```

## Complexity

★ **Running time**

$O(m \times \log n)$

★ **Worst case**

– list  $L = (\mathbf{a}^{m-1}\mathbf{b}, \mathbf{a}^{m-1}\mathbf{c}, \mathbf{a}^{m-1}\mathbf{d}, \dots)$

– string  $x = \mathbf{a}^m$

★ **Additional space**

constant

## Binary search tree

★ **Nodes**

$n + 1$  external nodes  $(-1, 0), (0, 1), (1, 2), \dots, (n - 1, n)$

$n$  internal nodes in the form  $(d, f)$  with  $d + 1 < f$

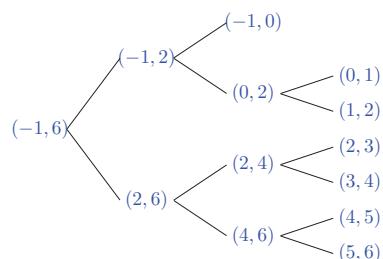
children of  $(d, f)$ :  $(d, \lfloor (d + f)/2 \rfloor)$  and  $(\lfloor (d + f)/2 \rfloor, f)$

root:  $(-1, n)$

★ **Size**

$2n + 1$  nodes for a list of  $n$  strings

★ **Example for  $n = 6$**



## Search using LCP's

- ★ **Aim**

reduce the running time to  $O(m + \log n)$

- ★ **LCP**, *longest common prefix*

$lcp(L_d, L_f)$  known for any pair  $(d, f)$  considered in the binary search

- ★ **Additional space**:  $O(n)$  integers for the  $2n+1$  LCP's associated with nodes of the binary search tree

- ★ **Algorithm** based on properties arising in three cases (plus symmetric cases)

- ★ **Variables**

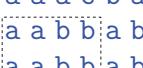
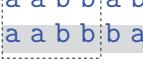
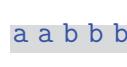
$ld = |lcp(x, L_d)|$ ,  $lf = |lcp(x, L_f)|$ ,  $i = \lfloor (d + f)/2 \rfloor$   
maintained during execution

## Case one

- ★ **Hypotheses**

$L_d < x < L_f$  and  $ld \leq |lcp(L_i, L_f)| < lf$

- ★ **Example**

$L_d$		$x$	
	a a a c b a		
$L_i$		$x$	
	a a b b a b b		
$L_f$		$x$	

- ★ **Conclusion**

$L_i < x < L_f$  and  $|lcp(x, L_i)| = |lcp(L_i, L_f)|$

## Case two

- ★ **Hypotheses**

$L_d < x < L_f$  and  $ld \leq lf < |lcp(L_i, L_f)|$

- ★ **Example**

$L_d$	a a a c a	x a a b a c b
	a a a c b a	
$L_i$	a a b b a b a	
	a a b b a b b	
$L_f$	a a b b b a b	x a a b a c b

- ★ **Conclusion**

$L_d < x < L_i$  and  $|lcp(x, L_i)| = |lcp(x, L_f)|$

## Case three

- ★ **Hypotheses**

$L_d < x < L_f$  and  $ld \leq lf = |lcp(L_i, L_f)|$

- ★ **Example**

$L_d$	a a a c a	x a a b b a b
	a a a c b a	
$L_i$	a a b b a b a	
	a a b b a b b	
$L_f$	a a b b b a b	x a a b b a b

- ★ **Conclusion**

compare  $x$  and  $L_i$  from position  $lf$

## Improved searching algorithm

```

SEARCH( $L, n, Lcp, x, m$ )
1  ( $d, ld$ )  $\leftarrow (-1, 0)$ 
2  ( $f, lf$ )  $\leftarrow (n, 0)$ 
3  while  $d + 1 < f$  do            $\triangleright$  Invariant :  $L_d < x < L_f$ 
4     $i \leftarrow \lfloor (d + f)/2 \rfloor$ 
5    if  $ld \leq Lcp(i, f) < lf$  then
6      ( $d, ld$ )  $\leftarrow (i, Lcp(i, f))$ 
7    elseif  $ld \leq lf < Lcp(i, f)$  then
8       $f \leftarrow i$ 
9    elseif  $lf \leq Lcp(d, i) < ld$  then
10   ( $f, lf$ )  $\leftarrow (i, Lcp(d, i))$ 
11   elseif  $lf < ld < Lcp(d, i)$  then
12      $d \leftarrow i$ 
13   else  $\ell \leftarrow \max\{ld, lf\}$ 
14    $\ell \leftarrow \ell + |lcp(x[\ell..m - 1], L_i[\ell..|L_i| - 1])|$ 
15   if  $\ell = m$  and  $\ell = |L_i|$  then
16     return  $i$ 
17   elseif  $(\ell = |L_i|)$  or  $(\ell \neq m \text{ and } L_i[\ell] < x[\ell])$  then
18     ( $d, ld$ )  $\leftarrow (i, \ell)$ 
19   else ( $f, lf$ )  $\leftarrow (i, \ell)$ 
20 return  $(d, f)$ 

```

## Complexity

**Proposition 1** Algorithm SEARCH finds a string  $x$  of length  $m$  in a sorted list of  $n$  strings in time  $O(m + \log n)$ .

It makes no more than  $m + \lceil \log(n + 1) \rceil$  comparisons of letters.

It requires  $O(n)$  extra space.

### Sketch of the proof

Number of letter comparisons:

- ★ each positive comparison strictly increases  $\ell$ , yielding no more than  $m$  such comparisons.
- ★ each negative comparison leads to divide by two the value of  $f - d$ , producing no more than  $\lceil \log(n + 1) \rceil$  such comparisons.

LCP can be implemented to run in constant time after preprocessing.

## Interval

```

1  ▷ next line replaces line 15 of SEARCH
2  if  $\ell = m$  then
3      ▷ next lines replace line 16 of SEARCH
4       $e \leftarrow i$ 
5      while  $d + 1 < e$  do
6           $j \leftarrow \lfloor (d + e)/2 \rfloor$ 
7          if  $Lcp(j, e) < m$  then
8               $d \leftarrow j$ 
9          else  $e \leftarrow j$ 
10         if  $Lcp(d, e) \geq m$  then
11              $d \leftarrow \max\{d - 1, -1\}$ 
12          $e \leftarrow i$ 
13         while  $e + 1 < f$  do
14              $j \leftarrow \lfloor (e + f)/2 \rfloor$ 
15             if  $Lcp(e, j) < m$  then
16                  $f \leftarrow j$ 
17             else  $e \leftarrow j$ 
18             if  $Lcp(e, f) \geq m$  then
19                  $f \leftarrow \min\{f + 1, n\}$ 
20         return  $(d, f)$ 

```

## Preprocessing the list

Let  $\|L\| = \sum_{i=0}^{n-1} |L_i|$ .

★ **Sorting**

repetitive application of bin sorting: time  $O(\|L\|)$

★ **Computing LCP's** of  $L_{f-1}$  and  $L_f$ ,  $0 \leq f \leq n$   
straight algorithm: time  $O(\|L\|)$

★ **Computing other LCP's**

based on next lemma

**Lemma 1** Let  $L_0 \leq L_1 \leq \dots \leq L_{n-1}$ . Let  $d, i$  and  $f$ ,  $-1 < d < i < f < n$ . Then  $|lcp(L_d, L_f)| = \min\{|lcp(L_d, L_i)|, |lcp(L_i, L_f)|\}$ .

**Proposition 2** Preprocessing  $L$ , sorting and computing LCP's, takes  $O(\|L\|)$  time.