# XML and Databases

Lecture 13 Fast Substring Search

Sebastian Maneth NICTA and UNSW

CSE@UNSW - Semester 1, 2010

Recall the **contai ns-predicate** of XPath:

```
//book/abstract[contains(., "fix")]
```

For instance the abstract node:

<book>. .

```
<abstract>This article dicusses the advantages of
suffix arrays, for the purpose of substring search ...
</abstract>...
```

</book>

will be returned, because it contains the substring "fix" because it appears in the word "suffix" mentioned in the abstract text.

#### Question

Given a very large text, how do you search for

- $\rightarrow$  All occurrences of a given keyword?
- $\rightarrow$  All occurrences of a given substring?
- $\rightarrow$  Count them (can be done faster?)

#### Question

Given a very large text, how do you search for

- $\rightarrow$  All occurrences of a given keyword?
- $\rightarrow$  All occurrences of a given substring?
- $\rightarrow$  Count them (can be done faster?)

What we know so far:

 $\rightarrow$  can use KMP-algorithm.

for a text of length n, it only takes O(n) time to locate all occurrences of the substring.

→ in a database, that is \*way\* to slow!!
How do you think Google indexes text for fast search??

#### Question

Given a very large text, how do you search for

- $\rightarrow$  All occurrences of a given keyword?
- $\rightarrow$  All occurrences of a given substring?
- $\rightarrow$  Count them (can be done faster?)

We want search time to be independent of the size n of the text, but should only depend on the length of the keyword.

We are allowed to preprocess the string in linear time ("indexing").

#### Question

Given a very large text, how do you search for

- $\rightarrow$  All occurrences of a given keyword?
- $\rightarrow$  All occurrences of a given substring?
- $\rightarrow$  Count them (can be done faster?)

Idea 1 --we search for exact WORDS, not substrings—



Given a keyword  $a_1a_2...a_m$  of length m,

How much time required to locate all occurrences of the keyword?

Easy: keep start rows of strings that "start with a\_1" (for any letter), and within those rows, again those that "continue with letter a\_2" (for all letters) Etc. (this is a tree of height=length of longest word, and branching=# different letters)

Idea 1 --we search for exact WORDS, not substrings—



Given a keyword  $a_1a_2...a_m$  of length m,

How much time required to locate all occurrences of the keyword?

```
→ only time O(m)! ☺
```

Problems (1) indexing time?! (2) how to do substring search??

Idea 1 --we search for exact WORDS, not substrings—



Given the text of length n, how many **substrings** are there?

 $\rightarrow$  (begin position, end position)

Quadratically many! That is, O(n^2). Thus, it is **impossible** in linear time to list all these substrings and put them into a (sorted) dictionary!



Idea 1 --we search for exact WORDS, not substrings—



Idea comes from compression. bzi p2 is based on the Burrows-Wheeler Transform!



Idea comes from compression.

bzi p2 is based on the Burrows-Wheeler Transform!

1) Add an end-marker "\$" to the end of the text 2) End-marker \$ is smallest in ordering: '\$' < 'a' < 'b' < 'c' < ..... < 'z' < 'A' < .... Compute all cyclic shifts of text 3) Sort them lexicographically 4) Burrows-Wheeler Transform of text T banana\$ \$banana \$banana a\$banan Question a\$banan ana\$ban Why do you think is the BWT good na\$bana anana\$b for compression? ana\$ban banana\$ sort nana\$ba na\$bana anana\$b nana\$ba

Idea comes from compression.

bzi p2 is based on the Burrows-Wheeler Transform!

1) Add an end-marker "\$" to the end of the text 2) End-marker \$ is smallest in ordering: '\$' < 'a' < 'b' < 'c' < ..... < 'z' < 'A' < .... Compute all cyclic shifts of text 3) 4) Sort them lexicographically Burrows-Wheeler Transform of text T banana\$ \$banana First row: only tells us \$banana a\$banan how many substrings a\$banan ana\$ban  $\rightarrow$  start with "a" (3) na\$bana anana\$b  $\rightarrow$  how many start with "b" (1) ana\$ban banana\$ sort etc. nana\$ba na\$bana Same for any text with these letters! anana\$b nana\$ba We canNOT reconstruct T from row 1!

Idea comes from compression. bzi p2 is based on the Burrows-Wheeler Transform!











#### **BWT: Better Decompression**

→ In a real implementation we may NOT construct all cyclic shifts and sort... (because that takes quadratic time!!)
 → Same for decompression. May not do it the naïve way!



Retrieving T: start from end marker, read backwards (by applying LF)

Here comes the *magic*: we are now able to count the number of occurrences of a substring of length m, only in time O(m log S)!

S = size of alphabet

This is what makes fast keyword Search a la Google possible!

Search time is INDEPENDENT of the size of the text!!

Here comes the *magic*: we are now able to count the number of occurrences of a substring of length m, only in time O(m log S)!

```
banana$
              $banana
                            C $ a b n
$banana
              a$banan
                              0 1 4 5
a$banan
              ana$ban
na$bana
              anana$b
ana$ban
              banana<sup>$</sup>
                            LF-mapping
                            LF(i)=C[L[i]] + rank_{L[i]}(L,i)
nana$ba
              na$bana
anana$b
              nana$ba
                                                     O(log S) time
                                                     using wavelet tree
```

Backward search for Pattern P[1]..P[m]

Initial range: [sp,ep] with sp=C[P[m]]+1 and ep=C[P[m]+1] Then [s,e] with s = C[P[i]] + rank<sub>L[i]</sub>(L, sp-1) + 1 e = C[P[i]] + rank<sub>L[i]</sub>(L, ep) S = size of

alphabet

Here comes the *magic*: we are now able to count the number of occurrences of a substring of length m, only in time O(m log S)!



#### Backward search for Pattern P[1]..P[m]

→ Initial range: [sp,ep] with sp=C[P[m]]+1 and ep=C[P[m]+1] Then [s,e] with  $s = C[P[i]] + rank_{L[i]}(L, sp-1) + 1$  $e = C[P[i]] + rank_{L[i]}(L, ep)$  S = size of

alphabet

Here comes the *magic*: we are now able to count the number of occurrences of a substring of length m, only in time O(m log S)!



Backward search for Pattern P[1]..P[m]

Then [s,e] with  $s = C[P[i]] + rank_{L[i]}(L, sp-1) + 1$  $e = C[P[i]] + rank_{L[i]}(L, ep)$  22

S = size ofalphabet

 $e = 5 + rank_n(L,4)$ 

= 5 + 2 = 7

Here comes the *magic*: we are now able to count the number of occurrences of a substring of length m, only in time O(m log S)!

-7

S = size of

alphabet

23

Here comes the *magic*: we are now able to count the number of occurrences of a substring of length m, only in time O(m log S)!

S = size of alphabet

Backward search for Pattern P[1]..P[m]

Counting: O(m log S) time

**Locating** If every  $I = log^{1+epsilon} n$  position is sampled then O(I log S) per occurrence, by backward traversal using LF.

#### **Real Performance**

/\* In order : IsContains, Timing of IsContains, GlobalCount, Timing of GlobalCount, CountContains, time of CountContains, time of Full Report Contains \*/

#### Sampling rate 64

"Bakst": 1, 0.038, 1, 0.004, 1, 0.04, 0.012, max\_mem = 61 "rumi nants": 1, 0.04, 22, 0.009, 19, 2.281, 1.588, max\_mem = 61 "morphine": 1, 0.026, 392, 0.009, 144, 29.924, 32.668, max\_mem = 61 "AUSTRALIA": 1, 0.028, 438, 0.009, 438, 4.616, 4.457, max\_mem = 61 "mol ecul e": 1, 0.051, 1472, 0.008, 966, 128.28, 122.014, max\_mem = 61 "brain": 1, 0.02, 2685, 0.005, 1493, 218.462, 215.196, max\_mem = 61 "brain": 1, 0.019, 6897, 0.005, 4690, 553.496, 548.009, max\_mem = 62 "bl ood": 1, 0.018, 10402, 0.005, 8534, 401.214, 399.674, max\_mem = 62 "from": 1, 0.016, 20859, 0.004, 12073, 1722.95, 1717.83, max\_mem = 62 "wi th": 1, 0.016, 63322, 0.004, 22974, 5084.14, 5083.77, max\_mem = 63 "in": 1, 0.001, 2932251, 0, 595716, 189299, 188377, max\_mem = 93 "\n": 1, 0.001, 9730750, 0.001, 5870474, 132780, 132241, max\_mem = 86

use

naïve CountContains/FullContains on naïve text: ca. **2700ms** 

#### **Real Performance**

/\* In order : IsContains, Timing of IsContains, GlobalCount, Timing of GlobalCount, CountContains, time of CountContains, time of Full Report Contains \*/

Sampling rate 5 "Bakst": 1, 0.038, 1, 0.005, 1, 0.049, 0.013, max mem = 100 "ruminants": 1, 0.038, 22, 0.01, 19, 0.156, 0.086, max\_mem = 100 "morphine": 1, 0.027, 392, 0.009, 144, 1.718, 1.357, max\_mem = 100 "AUSTRALIA": 1, 0.098, 438, 0.009, 438, 4.145, 3.942, max mem = 100 "molecule": 1, 0.029, 1472, 0.009, 966, 6.247, 5.853, max mem = 101 "brain": 1, 0.019, 2685, 0.006, 1493, 12.24, 11.588, max mem = 101 "human": 1, 0.018, 6897, 0.005, 4690, 25.403, 27.344, max\_mem = 101 "blood": 1, 0.026, 10402, 0.005, 8534, 77.175, 73.613, max mem = 101 "from": 1, 0.016, 20859, 0.003, 12073, 84.012, 78.663, max\_mem = 101 "with": 1, 0.015, 63332, 0.004, 22974, 242.834, 235.043, max\_mem = 102 <u>" in": 1, 0.012, 238638, 0.002, 42586, 1105.6, 1091.43, max\_mem = 103</u> "b": 1, 0, 411409, 0.001, 135307, 1779.27, 1762.62, max mem = 108 "g": 1, 0.001, 748326, 0, 320440, 3411.65, 3378.85, max\_mem = 119 ""a": 1, 0, 2932251, 0, 595716, 13183.4, 13173.4, max\_mem = 133 "\n": 1, 0.001, 9730750, 0.001, 5870474, 87770.9, 88230.4, max mem = 126

use

naïve CountContains/FullContains on naïve text: ca. 2700ms

#### **Construction Time**

XMark data 174 different element labels Max Depth: 14, Average Depth: 9.6

116MB XMark 6,074,297 nodes 559MB XMark 29,239,763 nodes

1GB XMark 58, 472, 941 nodes

Text:	7min 18s TOTAL= 9min 20s
Text:	38min 45s TOTAL= 53min 25s
Text:	1h 24min TOTAL= 1h 55min



#### ☺ Advertisement ☺

New course, will be first offered in Session 1 of 2011.

COMP9319 -- Web Data Compression and Search (PG, UOC: 6)

Contents

**Data Compression :** (a) Adaptive Coding, Information Theory (b) Text Compression (ZIP, GZIP, BZIP, etc) (c) Burrows-Wheeler Transform and Backward Search

(d) XML Compression

Search: (a) Indexing

(b) Pattern Matching and Regular Expression Search

- (c) Distributed Querying
- (d) Fast Index Construction
- (e) Implementation

If time allows: Streaming Algorithms, On-Line Data Analytics

The lecture materials will be complemented by projects and assignments.

# END Lecture 13 and of the course.

- $\rightarrow$  Thanks for your attention and hard work.
- $\rightarrow$  Hopefully you have enjoyed the lecture.
- $\rightarrow$  Good luck and all the best with

the exam on June 12<sup>th</sup>.