

COT 6936: Topics in Algorithms

Giri Narasimhan
 ECS 254A / EC 2443; Phone: x3748
 giri@cs.fiu.edu
http://www.cs.fiu.edu/~giri/teach/COT6936_S10.html
<https://online.cis.fiu.edu/portal/course/view.php?id=427>

2/25/10 COT 6936 1

The String Matching Problem

Pattern **P** → [] → Set of Locations **L**
 Text **T** → []

2/25/10 COT 6936 2

Approximate String Matching

Input: Text **T**, Pattern **P**

Questions:

- Does **P** occur in **T**?
- Find one/all occurrence of **P** in **T**.
- Count # of occurrences of **P** in **T**.
- Find longest substring of **P** in **T**.
- Find closest substring of **P** in **T**.
- Locate direct repeats of **P** in **T**.

Many more variants

2/25/10 COT 6936 3

String Matching Algorithms

- Find all occurrences of **P** in **T**.
 - Naïve Method
 - Rabin-Karp Method
 - FSA-based method
 - Knuth-Morris-Pratt algorithm
 - Boyer-Moore
 - Suffix Tree method
 - Shift-And method
 - Suffix Arrays
 - Methods based on Burrows-Wheeler transform

2/25/10 COT 6936 4

Naïve Strategy

ATAQAANANASPVANAGVERANANESISITALVDANANANANAS

PPPPANANAS
ANANAS
ANANAS
ANANAS

2/25/10 COT 6936 5

Finite State Automata

ANANAS

ATAQAANANASPVANAGVERANANESISITALVDANANANANAS

2/25/10 COT 6936 6

State Transition Diagram

		A	N	S	*
-	0	1	0	0	0
A	1	1	2	0	0
AN	2	3	0	0	0
ANA	3	1	4	0	0
ANAN	4	5	0	0	0
ANANA	5	1	4	6	0
ANANAS	6	1	0	0	0

2/25/10

COT 6936

7

Sliding Window Strategies

```

Initialize window on T;
While (window within T) do
  Scan: if (window = P) then report it;
  Shift: shift window to right (by ?? positions)
endwhile;
    
```

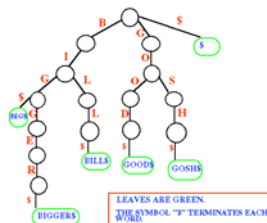
2/25/10

COT 6936

8

Tries

Storing:
 BIG
 BIGGER
 BILL
 GOOD
 GOSH



In this figure, the strings either start with B or G. Therefore, the root of the trie is connected to 3 edges called B, G and S.

1/2005

CAP5510/CGSS166 (Lec 4)

9

Suffix Tries & Compact Suffix Tries

Suffix Trie

Store all suffixes of
GOOGOLS

Compact Suffix Trie

2/25/10 COT 6936 10

Suffix Tries to Suffix Trees

Compact Suffix Trie

Suffix Tree

2/25/10 COT 6936 11

Suffix Trees

- **Linear-time** construction!
- String Matching, Substring matching, substring common to k of n strings
- All-pairs prefix-suffix problem
- Repeats & Tandem repeats
- Approximate string matching

1/2005 CAP5510/CGSS166 (Lec 4) 12

123456789012
 abracadabra\$

12 \$
 11 a\$
 8 abra\$
 1 abracadabra\$
 4 acadabra\$
 6 adabra\$
 9 bra\$
 2 bracadabra\$
 5 cadabra\$
 7 dabra\$
 10 ra\$
 3 racadabra\$

Search **abr** (len = m)

- Binary Search in suffix array
- $O(\log n)$ string comparisons
- Each comparison may involve comparing $O(m)$ characters
- $O(m + \log n)$ search
- $O(n \log n)$ space

2/25/10 COT 6936 13

Suffix Array & BWT Coding

12345678901 abracadabra

Space: $n \log n + O(n)$
 Search: $O(m + \log n)$

12345678901 abracadabra

Space: $n + O(A) + o(n)$
 Search: $O(m)$

- L is the BWT Code
- L is more compressible than the original input T

14

BWT Decoding: Step 1

Step 1: Compute F from L

F	L
12 \$ abracadabr <u>a</u>	12 \$ abracadabr <u>a</u>
11 a\$ abracadabr <u>a</u>	11 a\$ abracadabr <u>a</u>
8 abra\$ abracad <u>a</u>	8 abra\$ abracad <u>a</u>
1 abracadabra\$ <u>a</u>	1 abracadabra\$ <u>a</u>
4 acadabra\$ abr <u>a</u>	4 acadabra\$ abr <u>a</u>
6 adabra\$ abra <u>a</u>	6 adabra\$ abra <u>a</u>
9 bra\$ abracada <u>a</u>	9 bra\$ abracada <u>a</u>
2 bracadabra\$ a <u>a</u>	2 bracadabra\$ a <u>a</u>
5 cadabra\$ abra <u>a</u>	5 cadabra\$ abra <u>a</u>
7 dabra\$ abrac <u>a</u>	7 dabra\$ abrac <u>a</u>
10 ra\$ abracadab <u>a</u>	10 ra\$ abracadab <u>a</u>
3 racadabra\$ ab <u>a</u>	3 racadabra\$ ab <u>a</u>

15

Backward Search

Algorithm $\text{count}(P[1, p])$

```
(1)  $i \leftarrow p, c \leftarrow P[p], \text{First} \leftarrow 1, \text{Last} \leftarrow n;$ 
(2) while  $(\text{First} \leq \text{Last})$  and  $(i \geq 1)$  do
(3)    $c \leftarrow P[i];$ 
(4)    $\text{First} \leftarrow C[c] + \text{Occ}(c, \text{First} - 1) + 1;$ 
(5)    $\text{Last} \leftarrow C[c] + \text{Occ}(c, \text{Last});$ 
(6)    $i \leftarrow i - 1;$ 
(7) if  $(\text{Last} < \text{First})$  then return "no rows prefixed by  $P[1, p]$ " else return  $(\text{First}, \text{Last}).$ 
```

Fig. 3. Algorithm count for finding the set of rows prefixed by $P[1, p]$, and thus for counting the pattern occurrences $\text{occ} = \text{Last} - \text{First} + 1$. Recall that $C[c]$ is the number of text characters which are alphabetically smaller than c , and that $\text{Occ}(c, q)$ denotes the number of occurrences of character c in $T^{\text{bwt}}[1, q]$.

Backward Search

Need 2 mappings:

- $C: \Sigma \rightarrow [1, n]$ **EASY**
- $\text{Occ}: \Sigma \times [1, n] \rightarrow [1, n]$ **NON-TRIVIAL**

20

Rank Query

Need to answer rank queries:

$\text{Occ}(c, L, i)$: # of occurrences of c in $L[1..i]$

- Compute a bit sequence B_c
 - 1, if corresponding character is c
 - 0, otherwise
- Computing Occ is reduced to computing number of 1's in B_c until position i

21

Rank Query

Need to answer rank queries:

$Occ(1, B, i)$ in $O(1)$ time with the minimal amount of space.

This is sufficient because:

- $Occ(1, B, i) = i - Occ(0, B, i)$
- Can be easily generalized to non-binary case using **Wavelet Trees**

22

Wavelet Trees

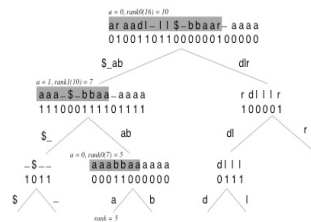


Fig. 5. A binary wavelet tree for the string $L = \text{"araad11l1\$bbaar_aaaa"}$, illustrating the solution of query $Occ(\text{"a"}, 15)$. Only the bit vectors are stored, the texts are shown for clarity.

23

Huffman-based Wavelet Trees

- Binary tree in the wavelet tree can be replaced with a **Huffman-based tree**, i.e., more frequent characters are placed at shallower leaves.

2/25/10

COT 6936

24

Rank query: binary sequences

Use a 2-level dictionary

- Chop into short sequences (of length $k = (\log n)/2$) and store solutions of ranks at regular intervals
- Store global table with answers for every possible short sequence. This table is of size $\sqrt{n} \times k$

2/25/10

COT 6936

25
