

# STRINGS

Motivating Problem Given text "string"  $T = t_1 t_2 \dots t_n \in \Sigma^*$

build a <sup>static</sup> data structure s.t.  
 given "pattern" string  $P = p_1 p_2 \dots p_m$ , ( $m \leq n$ )  
 decide whether  $P$  is a substring of  $T$   
 (if so, report one occurrence  
 or all occurrences  
 or count)

eg.  $T = 01011011010$   
 $P = 1011$

Non-DS version:  
 (Static string matching problem)

trivial alg'm  $O(mn)$  time  
 Knuth-Morris-Pratt '77  $O(n)$  time  
 (many other  $O(n)$  alg'ns...  
 eg. Karp-Rabin '87 in CS473)

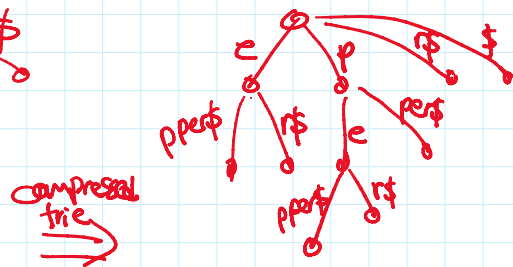
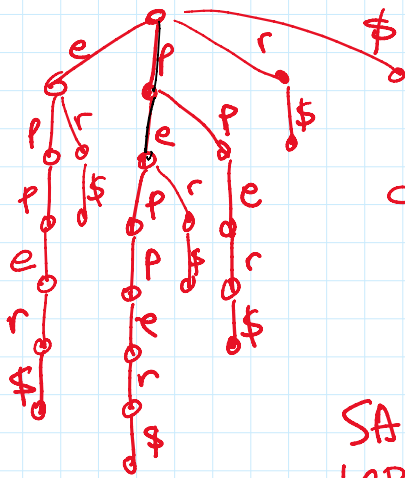
DS version?

## Suffix Tree

compressed trie for all suffixes of  $T$ .

eg.  $T = \text{"peppers\$"}$

- 1 peppers\$
- 2 eppers\$
- 3 pper\$
- 4 pers\$
- 5 ers\$
- 6 r\$
- 7 \$



compressed trie →

$O(n)$  size

$P = \text{"ep"}$   
 $\text{"pe"}$



query  $O(m)$   
 or  $O(m + \text{occ})$ .

SA: 2, 5, 1, 4, 3, 6, 7  
 LCP: 1, 0, 2, 1, 0, 0

$\Delta/\Delta \Delta$  or  $O(n \log n)$ .

Preproc time?? (algn for constructing suffix tree)

naive  $O(n^2)$

*Complicated* → Werner '73  $O(n)$  for  $|\Sigma| = O(1)$   
 $O(n \log n)$  for general  $\Sigma$ .

Farach-Colton '97  $O(n)$  for general  $\Sigma \subseteq [n]$   
(in word RAM)

## Simplification: Suffix Array

Just store sorted order of all the suffixes of  $T$

i.e.  $SA[i] =$  <sup>index of</sup>  $i^{\text{th}}$  suffix in sorted order

also store  $LCP[i] =$  length of longest common prefix  
between  $i^{\text{th}}$  &  $(i+1)^{\text{th}}$  suffix  
in sorted order.

Obs Suffix tree  $\iff$  suffix array + LCP

Pf:  $(\Leftarrow)$  repeated RMQ on LCP array.  $\square$

Algs to construct suffix array (+ LCP array):

naive  $O(n \log n \cdot n) = O(n^2 \log n)$

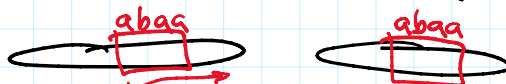
Gonnet et al. '92 / Manber-Myer '95  $O(n \log n)$

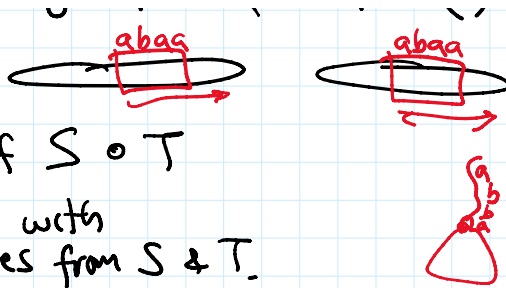
→ Kärkkäinen-Sanders '03  $O(n)$  *simple!*

⋮

Many many applications:

App 1 largest common substring of  $S$  &  $T$  in  $O(n)$  time





Compute suffix tree of  $S \circ T$   
 find deepest node with  
 (leaves from  $S$  &  $T$ .)

APP2 longest palindrome substring in  $S$  in  $O(n)$  time

compute suffix tree of  $S \circ S^R$  find hannahscat  
 for each  $i$ , find LCA of  $i^{\text{th}}$  suffix &  $(2n-i)^{\text{th}}$  suffix

a  
 v  
 c

### Algm for Suffix Array by Kärkkäinen & Sanders:

idea - recursion

assume  $\Sigma \subseteq [n]$ .

Fact sort  $[n] \times [n]$   
 or  $[n] \times [n] \times [n]$   
 in radix sort  
 in  $O(n)$  time

1. let  $T_1 = (\overset{\text{single char}}{t_1} t_2 t_3) (t_4 t_5 t_6) (t_7 t_8 t_9) \dots$   
 viewed as a string over  $\Sigma^3$  of length  $n/3$

can be mapped down to  $[n/3]$   
 by replacing with ranks.

$T_2 = (t_2 t_3 t_4) (t_5 t_6 t_7) (t_8 t_9 t_{10}) \dots$

$T_3 = (t_3 t_4 t_5) (t_6 t_7 t_8) (t_9 t_{10} t_{11}) \dots$

2. recurse in  $T_1 \circ T_2$

3. sort suffixes of  $T_3$  by radix sort over  $[n] \times [n]$  in

3. Sort suffixes of  $T_3$  by radix sort over  $[n] \times [n]$  in  $O(n)$  time.

how to compare

$$(t_{3i+3}, t_{3i+4}, t_{3i+5}, \dots) \stackrel{?}{<} (t_{3j+3}, t_{3j+4}, t_{3j+5}, \dots)$$

suffix in  $T_1$ 
suffix in  $T_1$

4. merge sorted list of suffixes of  $T_3$  with suffixes of  $T_1 \circ T_2$  in  $O(n)$  time

how to compare

$$(t_{3i+3}, t_{3i+4}, t_{3i+5}, \dots) \stackrel{?}{<} (t_{3j+1}, t_{3j+2}, t_{3j+3}, \dots)$$

suffix in  $T_1$ 
suffix in  $T_2$

$$(t_{3i+3}, t_{3i+4}, t_{3i+5}, \dots) < (t_{3j+2}, t_{3j+3}, t_{3j+4}, \dots)$$

suffix in  $T_2$ 
suffix in  $T_1$

$$\Rightarrow T(n) = T\left(\frac{2n}{3}\right) + O(n)$$

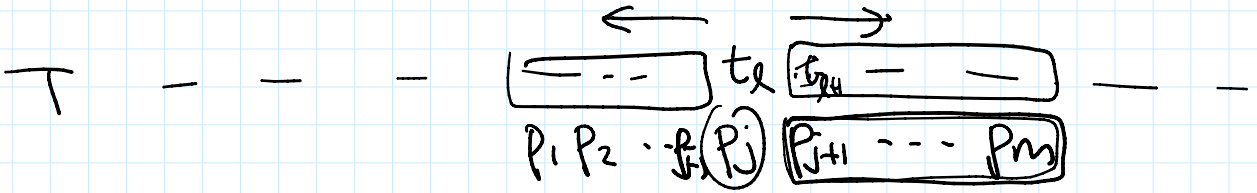
$$\Rightarrow O\left(n + \frac{2n}{3} + \frac{4}{9}n + \dots\right)$$

$$= \boxed{O(n)}$$

(can also compute LCP array ...)

More App 3: given query pattern  $P = p_1 \dots p_m$ ,  
decide whether  $\exists$  match with  $\leq 1$  error

e.g.  $T = \text{"this is a string"}$   
 $P = \text{"thin"}$



for  $j=1, \dots, m$ ,

search suffix tree of  $T$  for  $P_{j+1} \dots P_m$  ←

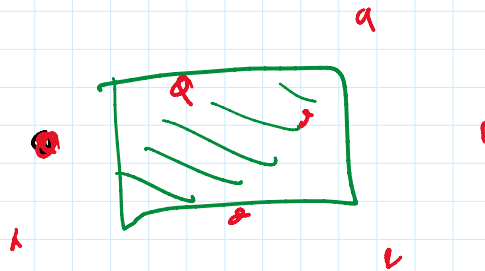
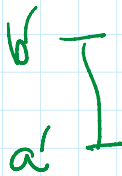
⇒ interval  $[a, b]$  in  $SA_T$

search suffix tree of  $T^R$  for  $P_j \dots P_1$  ←

⇒ interval  $[a', b']$  in  $SA_{T^R}$

find  $l$  s.t.  $\begin{cases} \text{position of } l+1 \text{ in } SA_T \in [a, b] \\ \text{position of } n-l+1 \text{ in } SA_{T^R} \in [a', b'] \end{cases}$

⇒ find point  $(x(l), y(l))$  in 2D rect.



⇒ 2D orth. range search

⇓

2D range tree

$O(n)$  space

query time  $O(\cancel{m^2} + m \log n)$   
 $m$   
 (with suffix links)