

Lecture 7 6.851

André Schult's Notes

STRING DS

- Notation :
- Σ , finite Alphabet
 - $T \in \Sigma^*$ text
 - $P \in \Sigma^*$ search pattern
 - $T[i:j]$ substring of T from character $i \dots j$
 - $T[i:]$ suffix starting at i
 - concatenation

STRING MATCHING

Input: text T , pattern P
Q: Find all^{*} occurrences of P in T

* (also some, first, first k...)

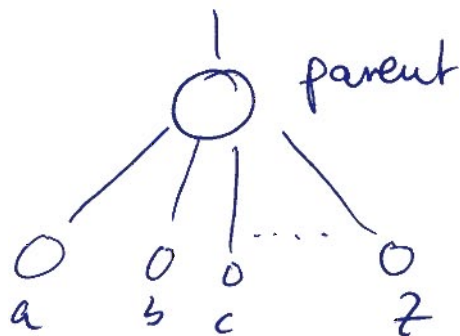
- 1 shot approach takes $O(n)$ time

(KNUTH, MORRIS, PRATT 77,
Boyer Moore 77,
Rabin, Karp 87)

How to "store" strings?

TRIE : • tree with fan-out $\leq |\Sigma|$

- edges parent \rightarrow child are labeled with the letters of Σ'



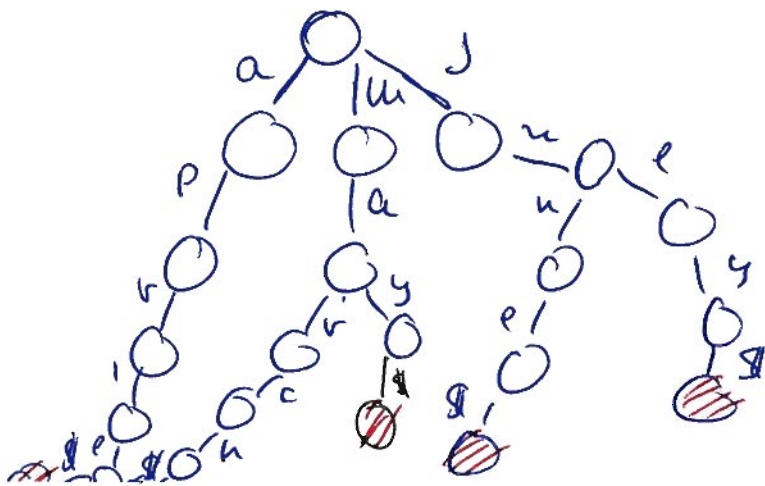
$\Sigma' = \{a, b, \dots, z\}$

- common technique: terminate strings with "\$"

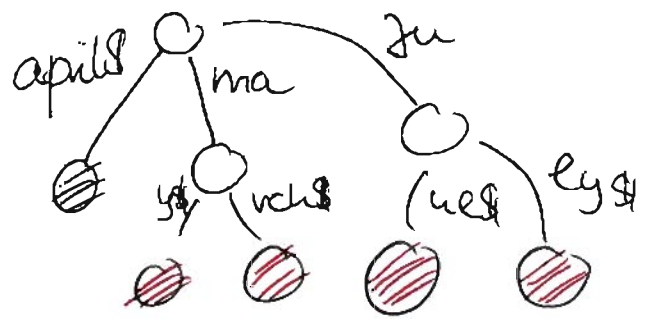
(to distinguish if a path represents a prefix or word)

- Root \rightarrow leaf paths corresponds to strings, stored in the trie

- example : { march, april, may, june, july }



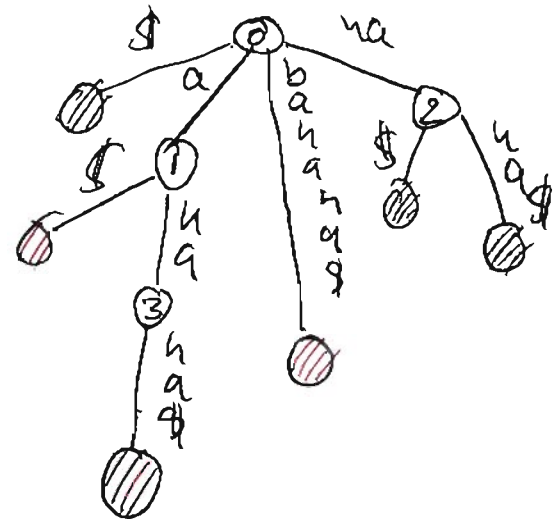
Compressed trie:
 (contract non-branching paths to single edges)



Suffix tree

: Compressed trie of all $|T|+1$ suffixes of T

• example: banana\$
 0 1 2 3 4 5 6



• ~~length of edge~~ leaf depth
 l. depth = length of the prefix stored at node v

• $|T|+1$ leaves representing the suffixes

• we store the edge labels $T[i:j]$ by storing i, j

↳ $O(|T|)$ space [words] not bits

Applications - Suffix Trees

- Find all occurrences of $P \in T$
(substring = prefix of a suffix) $O(|P|+k)$
- Count all occurrences of $P \in T$
(~~add~~ store subtree size in the interior nodes)
- longest repeated substring $\in O(|T|)$
(branching node with max. letter depth)
- multiple documents
(concatenate texts with $\$, \$_1, \$_2, \$_3, \dots, \$_k$)
- longest common substring $O(|T|)$
- max letter depth with ≥ 2 distinct $\$$'s below

Suffix arrays

- sort suffixes in T and store their indices (lexicographic order)

• Example :

6	\$	
5	a	\$
3	ana	\$
1	anana	\$
0	banana	\$
4	na	\$
2	nana	\$

LCP-Array

0
1
3
0
0
2

- Searchable in $O(|T| \log |T|)$ via binary search (~~speed up: homework~~)
- LCP Array : we store in cell i the length of the ^{common prefix} i -th and $(i+1)$ th entry of the suffix array
- Suffix + LCP Array : searchable in $O(|T| + \log |T|)$ with RMQ-DS (Homework)
- ~~LCP Array stores letter depth~~
- LCP Array can be constructed in $O(T)$ time or with the construction we discuss later for the suffix-Array [Kasai et al.]

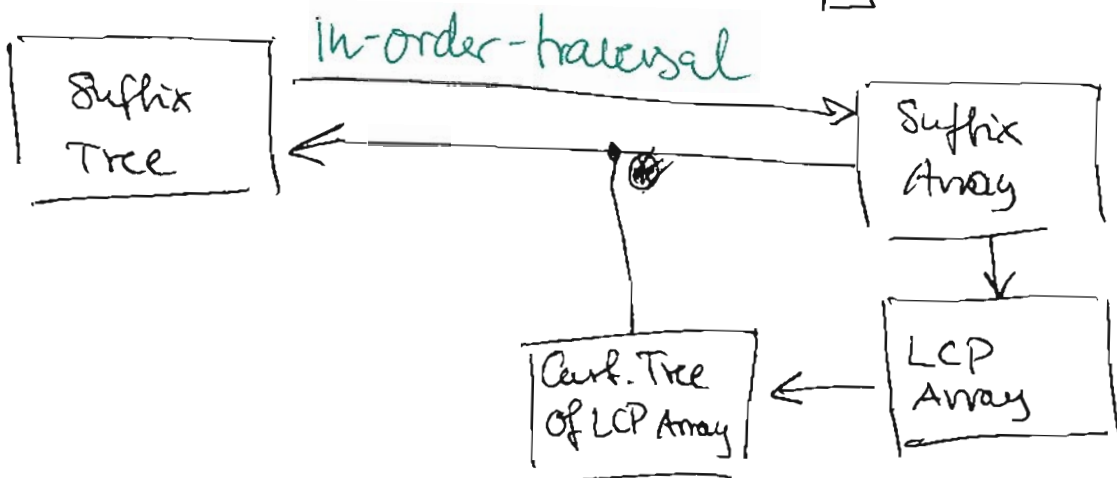
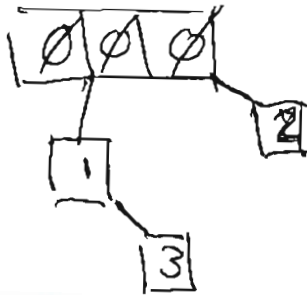
Suffix Tree \leftrightarrow Suffix Array

- ST & SA can be transformed into each other in linear time

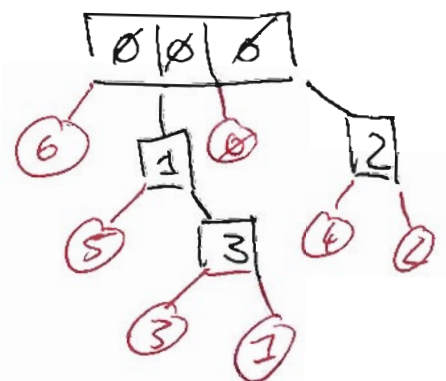
Cartesian-Tree (of LCP Array)

- Put all mins at the root
- recurse in remaining array pieces

- example:
(banana)



⊗ fill in suffixes (example)



- LCP's give letter depth
= length of edge labels

Building Suffix Array in $O(n)$ DC3 Alg.

$O(n) + \text{sort}(|\Sigma|)$

[Kärkkäinen Sanders - 03] ^{Burkhard} inspired by
 [Farach 97, Farach-Colton, Ferragina, Muthukrishnan 2000]

① Sort Σ in $|\Sigma| \log |\Sigma|$ time
 (remark later sorting uses Radix Sort $(O(n))$)

② Replace every letter in T by its rank in Σ

③ ^{Transform} Split T into 3 parts
 $T_0 = \langle (T[3i], T[3i+1], T[3i+2]), i=0,1,\dots \rangle$
 $T_1 = \langle (T[3i+1], T[3i+2], T[3i+3]), i=0,1,\dots \rangle$
 $T_2 = \langle (T[3i+2], T[3i+3], T[3i+4]), i=0,1,\dots \rangle$

④ Recurse on $T_0 \circ T_1$ (concatenated)
 (Distance cover)

↳ we get the suffixes in T sorted

⑤ Radix-Sort the suffixes in T_2 , use
 $T_2[i:] = T[3i+2:] = T[3i+2] \circ T[3i+3]$
 $= T[3i+2] \circ T_0[i+1]$

Toy example :

T = bananas \$\$\$
 0 1 2 3 4 5 6 7 8

$\{ \emptyset, a, b, c, d, s \}$

$T_0 = [ban] [ana] [s$$$]$
 c A F

$T_1 = [ana] [ues]$
 A D

$T_2 = [nan] [as$$$]$
 E B

$\emptyset < A < B < C < D < E < F$

$\tilde{T} = CAFAD$
 0 1 2 3 4 5

Recurse

Suffix Array SA(\tilde{T})

5	\$
1	AD\$
3	AFAD\$
0	CAFAD\$
4	D\$
2	FAD\$



i pos in SA

0	4.
1	3.
2	6.

Rank T_0

i pos in SA

0	2.
1	5.
2	1.

Rank T_1

To sort T_2 , write T_2 in Σ^* Suffixes (\tilde{T})

$$EB = [nan][as$$$] = (n, [ana][s$$$]) = (n, AF)$$

$$B = [as$$$] = (a, [s$$$]) = (a, F)$$

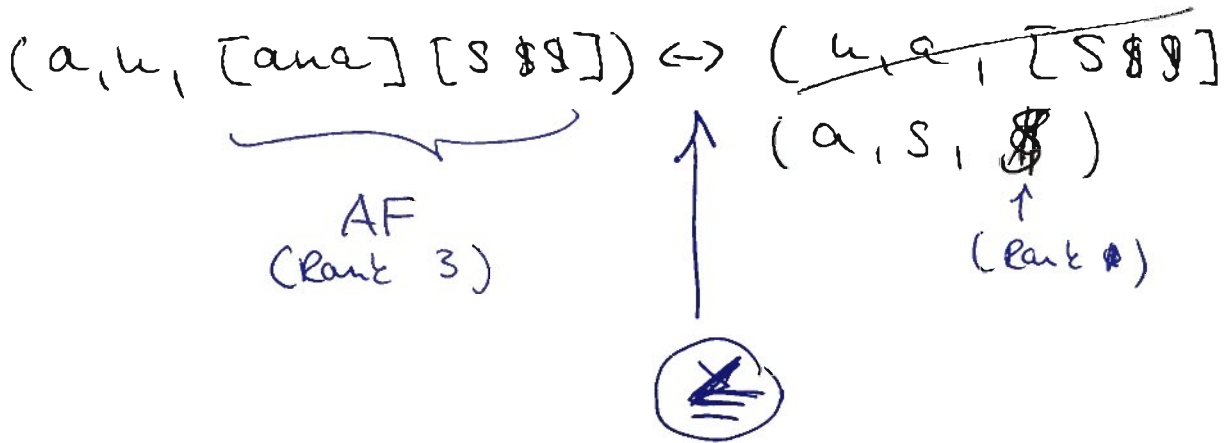
Radix Sort

SA for T_2 (ignore $\$$)

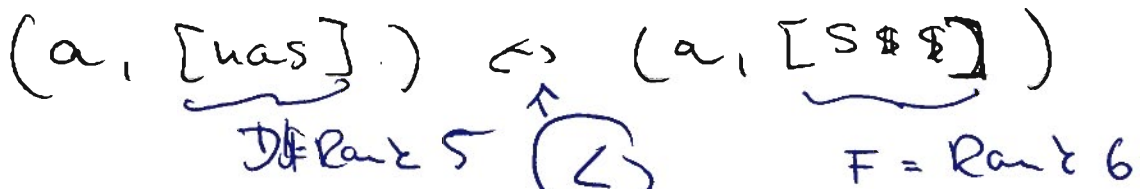
SA(T_2)	B	= (u, AF)
	EB	= (a, F)

Merging (ignore $\$$ entries)

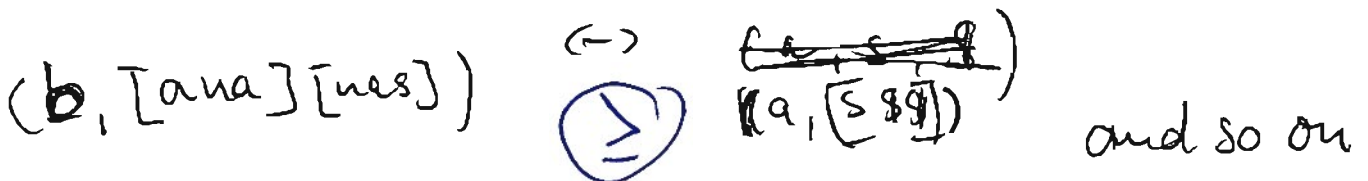
1. ~~✗~~ SA(\tilde{T}) [1] \leftrightarrow SA(T_2) [1]
 T_1 [0] \leftrightarrow T_2 [1] (Case B)



2. SA(\tilde{T}) [2] \leftrightarrow SA(T_2) [1]
 T_0 [1] \leftrightarrow T_2 [1] (Case A)

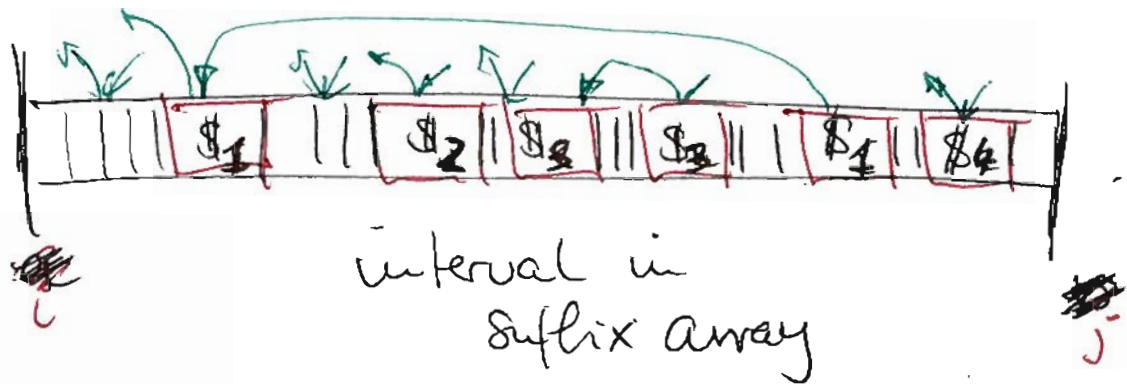


3. SA(\tilde{T}) [3] \leftrightarrow SA(T_2) [1]
 T_0 [0] \leftrightarrow T_2 [1] (Case A)



Document Retrieval [Muthukrishnan 2002]

- extract k distinct documents in $O(|P| + k)$ time interval $[i, j]$
- Patterns refer to ~~subtrees~~ in the suffix ~~tree~~ array
= (subtrees in suffix-tree)



Points: Position of the ~~text~~ $\$$ with the same index previous

↳ Stored in Array PA

- We want to find the first occurrences of $\$x$ only (for all x)
- We use a Range Minimum DS for the points

• find position of minimum stored at index k in the pointer array

→ if $PA[k] < i$: output k ,
recurse on $[i, k-1], [k+1, j]$
(with same i to left)