



Lecture:

SUFFIX ANCESTORY AND SUFFIX TREES

Unit:

9 - SUFFIX PROCESSING + COMPUTATIONAL LINGUISTICS

Instructor:

Wiesław

9.2.1 SUFFIX TREE AND APPLICATIONS

INFORMATION RETRIEVAL → pronounced "TRY"

 TREE OF ALL POSSIBLE

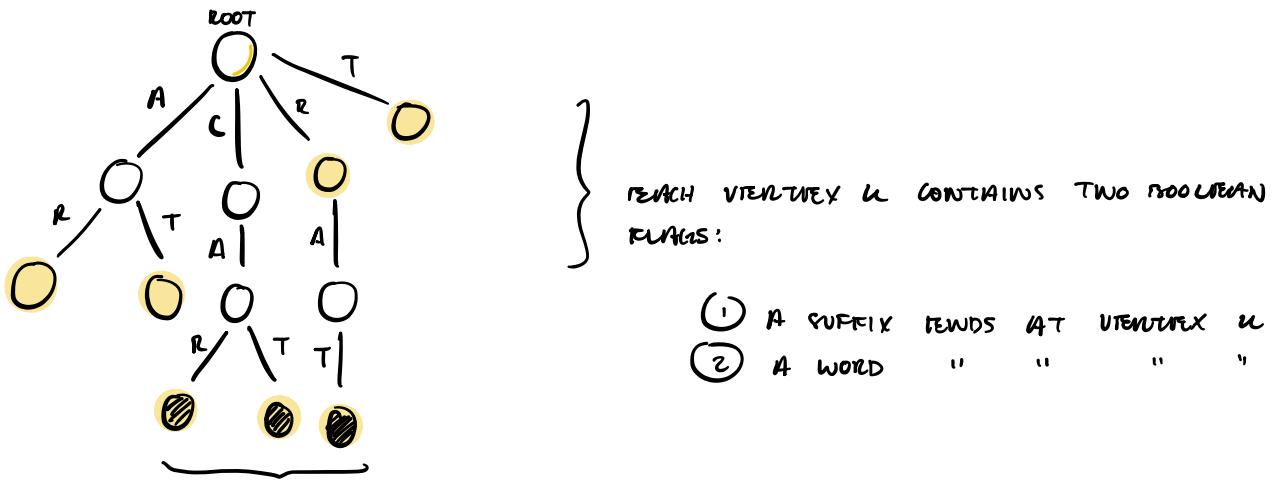
SUFFIXES OF A SET OF STRINGS S → DICTATE: REPRESENTMENTS IN CHARACTER
 UTENRDX: " " " SUFFIX REPRESENTED
 BY ITS PATH LABEL

$S = \{ \text{"CAT"}, \text{"CAT"}, \text{"RAT"} \} \rightarrow \text{DICTIONARY}$

SUFFIXES OF $S = \{ \text{'R}', \text{'AT'}, \text{'CAT'}, \text{'T'}, \text{'AT'}, \text{'CAT'}, \text{'RAT'} \}$

SUFFixed SUFFIXES OF $S = \{ \text{'AT'}, \text{'AT'}, \text{'CAT'}, \text{'CAT'}, \text{'R'}, \text{'RAT'}, \text{'T'} \}$

SUFFIX TREE:



THIS PATHS APPPEAR
IN THE DICTIONARY

APPLICATIONS:

- EFFICIENT DATA STRUCTURE FOR DICTIONARY
- DETERMINE IF A QUERY/PATTERN P EXISTS IN THE DICTIONARY IN $O(|P|)$

TRAVERSAL FROM THE ROOT DOWNWARDS

9.2.2 SUFFIX TENSE

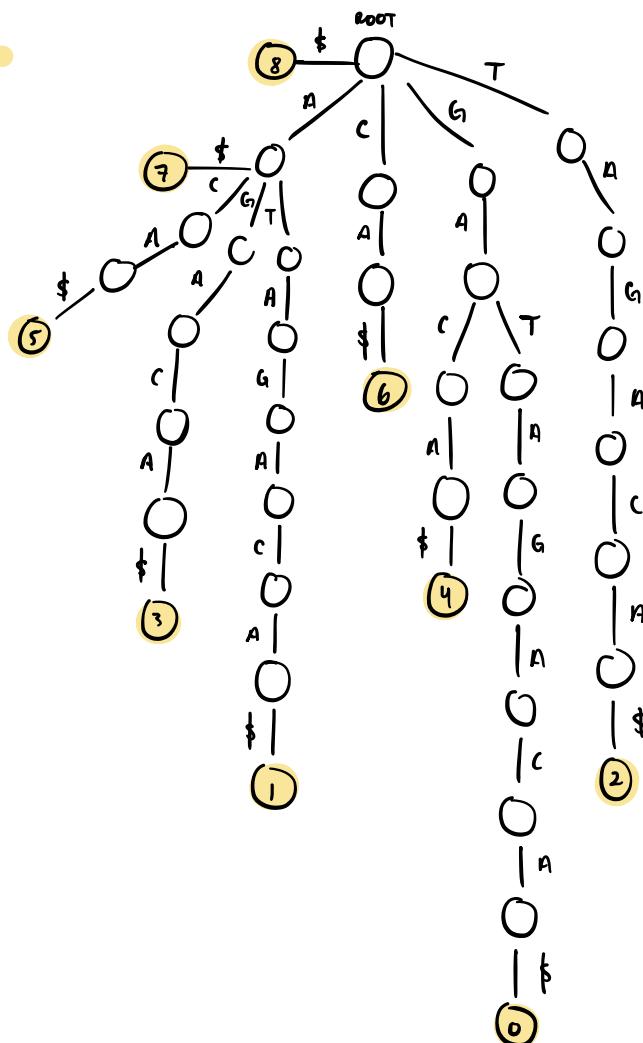
- ONE OF THE PROBLEMS OF A SUFFIX TREE IS THERE NUMBER OF REPETITIVE VERTICES FOR A GIVEN STRING

{ → SUFFIX THREE SOLUTES
THIS PROBLEM

→ T = "GATAGACCA \$"

↳ SPECIAL TERMINATING CHARACTER

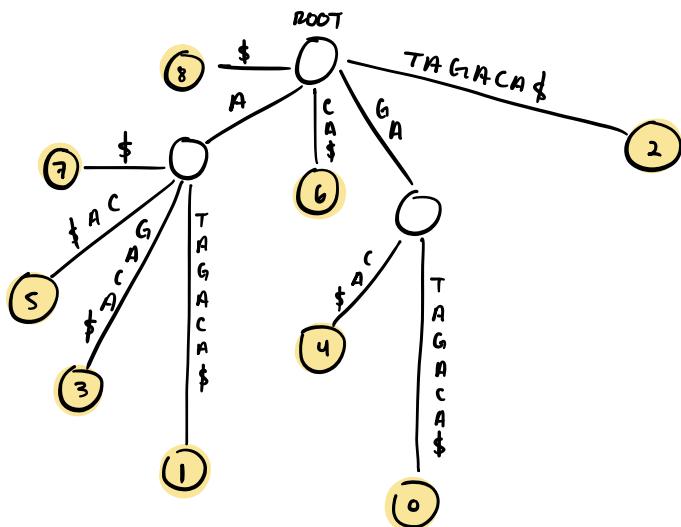
i	SUFFIX
0	GATAGACAS\$
1	ATAGACAS\$
2	TAGACAS\$
3	AGACAS\$
4	GACAS\$
5	ACAS\$
6	CAS\$
7	A\$
8	\$



SUFFIX TRUE OF T

- SUFFIX TRIE → SUFFIX TRIE WITHIN THE MANAGED VERTICES WITH ONLY ONE CHILD

- REGRE UNSTABLE CAN HAVE MORE THAN ONE CHARACTER
- AT MOST 2^n VERTICES



SUFFIX TREE OF T

9.2.3 APPLICATIONS OF SUFFIX TREE (ASSUMING SUFFIX TREE IS ALREADY BUILT)

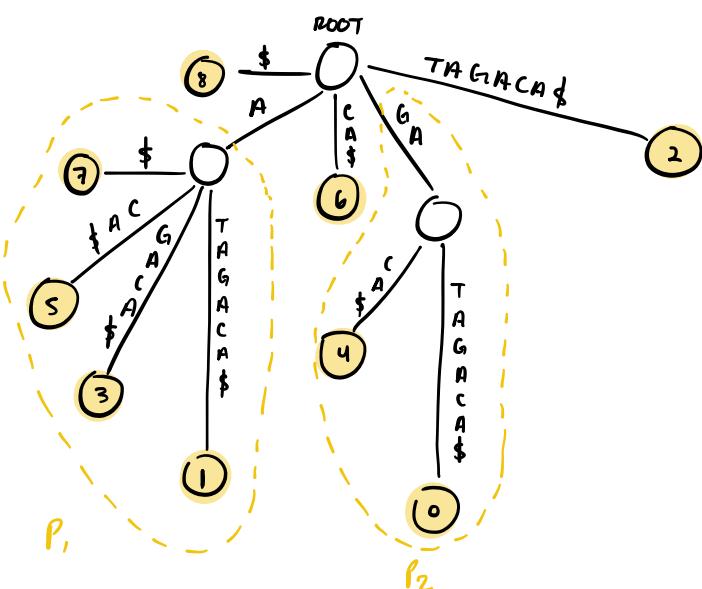
- [String Matching] in $O(m + \alpha c)$

m : LENGTH OF PATTERN P

- A MATCH IS A COMMON PREFIX

BETWEEN PATTERN STRING P AND SOME SUFFIXES OF STRING T

- FIND ALL OCCURRENCES OF STRING P IN STRING T



$T = "GATAGACA\$"$

0 1 2 3 4 5 6 7 8

$P_1 = 'A'$ \rightarrow OCCURRENCES : {7, 5, 3, 1}

$P_2 = "GA"$ \rightarrow " : {4, 0}

$P_3 = 'T'$ \rightarrow " : {2}

$P_4 = 'Z'$ \rightarrow " : { }

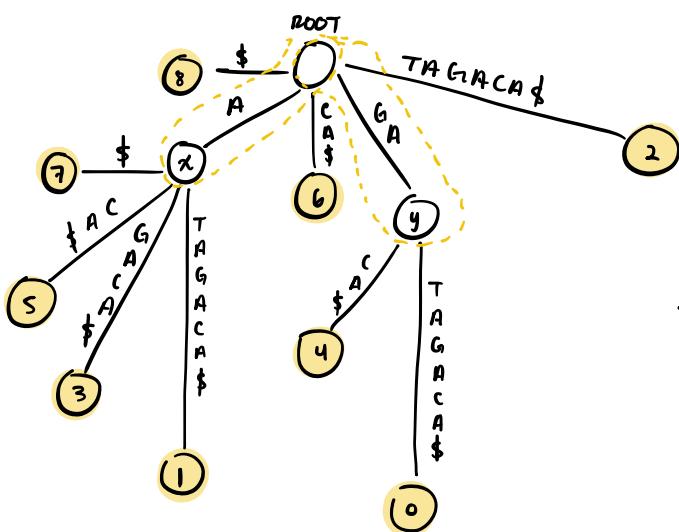
- [FINDING TIME LOWGEST FREQUENTED SUBSTRINGS] IN $O(n)$

- Answer: LENGTH OF PATH LARGEST OF THAT DEEPEST INTERVAL
 VERTICES X IN THE SUFFIX TREE



 PREFIX FOUND IN $O(n)$ USING THE TRAVERSAL

- INTERNAL VARIETY IN SUFFIX TENSE IMPLIES THAT IT REPRESENTS MORE THAN ONE SUFFIXES OF T



→ INTERNAL VERTICES:

- 1) X : PATH LENGTH LENGTH = 1 : A
2) y : " " " = 2 : GA

"GA" IS THE LRS

- [FINDING THE LONGEST COMMON SUBSTRING] in $O(n)$

- LCS OF TWO OR MORE STRINGS
 - BUILD A SUFFIX TREE THAT COMBINES THE SUFFIX TREES OF TWO STRINGS T_1 AND T_2
 - ↪ USE USE A TERMINATING SYMBOL FOR EACH T_1 AND T_2 (\$, #)

2) IDENTIFY INTERMEDIATE VERTICES THAT HAVE VERTICES IN THEIR SUBTRACES WITH DIFFERENT / TERMINATING SYMBOLS

COMMON PREFIXES BY STRINGS T_1 AND T_2

- 3) ANSWER IS THE PATH LARGEST WITH THE SHORTEST INTERVAL LENGTH

$T_1 = "GATAGACA\$"$ $T_2 = "CATATA\#"$

i	SUFFIX
0	GATAGACA\$
1	ATAGACA\$
2	TAGACA\$
3	AGACA\$
4	GACA\$
5	ACA\$
6	CA\$
7	A\$
8	\$

$$A : 1, 3, 5, 7 \\ 10, 12$$

$$C : 6, 9$$

$$G : 0, 4$$

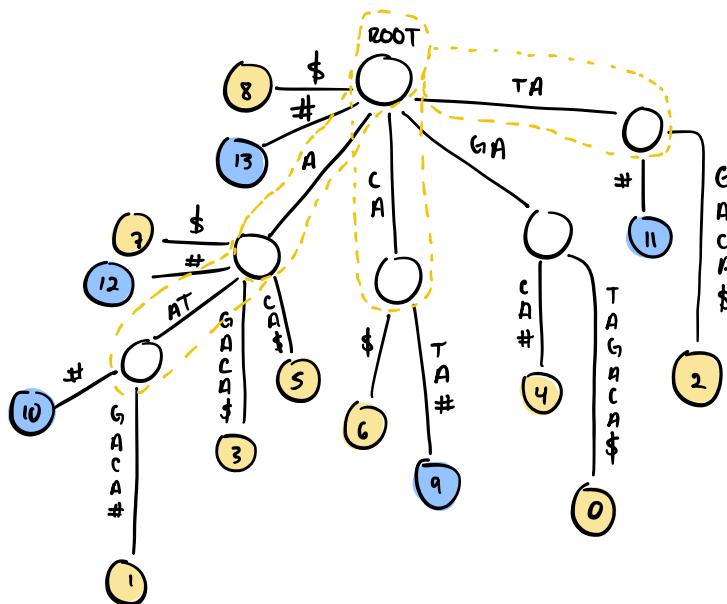
$$T : 2, 11$$

i	SUFFIX
9	CATA#
10	ATA#
11	TA#
12	A#
13	#

INTERNAL IDENTITIES TO CONSIDER:

- 1) PATH LENGTH $|A| = 1$
- 2) " " $|CA| = 2$
- 3) " " $|TA| = 2$
- 4) " " $|ATA| = 3$

ANSWER IS ATA



9.2.4 SUFFIX ARRAY

↳ MUCH SIMPLER TO BUILD THAN SUFFIX TREE $O(n \log n)$ CONST. TIME

→ SUFFIX ARRAY: INTEGRAL ARRAY THAT STORES A PERMUTATION OF N INDICES OF SUFFIXES

$$T = "GATAGACA\$"$$

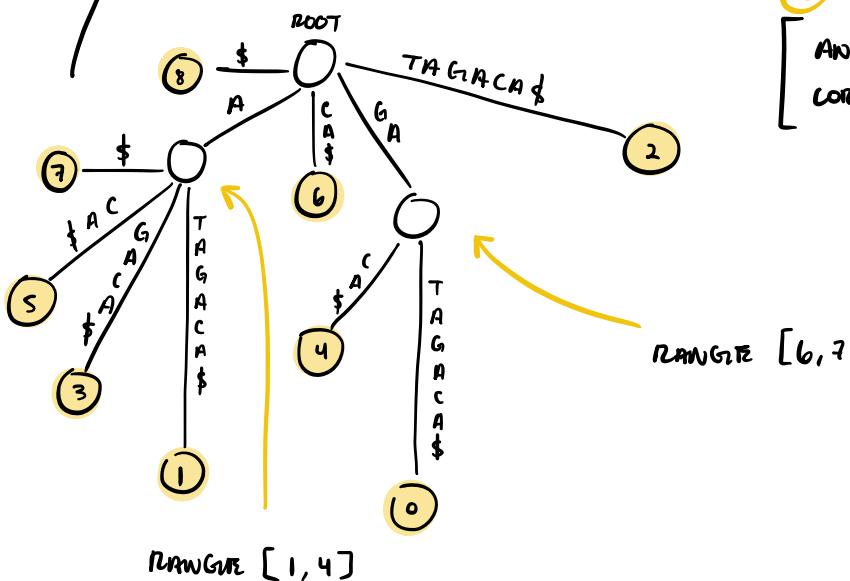
$$n = 9$$

$$SA = \{8, 7, 5, 3, 1, 6, 4, 0, 2\}$$

(1)

THE IN-PLACE TRAVERSAL OF THE SUFFIX TREE VISITS THE LEAVES IN SUFFIX ARRAY ORDER

i	SA[i]	SUFFIX
0	8	\$
1	7	A\$
2	5	ACA\$
3	3	AGACA\$
4	1	ATAGACA\$
5	6	CA\$
6	4	GA\$
7	0	GA\$
8	2	TAGACA\$



(2) AN INTERVAL INDEX IN ST CORRESPONDS TO A RANGE IN SA

COLLECTION OF SORTED SUFFIXES THAT SHARE A COMMON PREFIX

(3) A TERMINATING VERTICES (LEAF) IN ST CORRESPONDS TO AN INDIVIDUAL INDEX IN THE SA

9.2.5 APPLICATIONS OF SUFFIX ARRAY

- [STRING MATCHING] in $O(m \log n)$

m : LENGTH OF PATTERN P
 n : " " " STRING T

- PERFORM TWO $O(\log n)$ BINARY SEARCHES ON SORTED SUFFIXES AND UP TO $O(m)$ SUFFIX COMPARISONS

TO FIND SMALLEST AND LARGEST i SUCH THAT THE PREFIX OF $SA[i]$ MATCHES PATTERN P

- ALL OF THE SUFFIXES BETWEEN LOWER AND UPPER BOUND ARE OCCURRENCES OF P IN T