

Design & Analysis of Algorithms COMP 482 / ELEC 420



John Greiner

String Matching

Pattern: **CCATT**

Text: **ACTGCCATTCCTTAGGGCCATGTG**

- Brute force & variant
- Automata-like strategies
- Suffix trees

To do:
[CLRS] 32
Supplements
#5

Some Applications

- Text & programming languages
 - Spell-checking
 - Linguistic analysis
 - Tokenization
 - Virus scanning
 - Spam filtering
 - Database querying
- DNA sequence analysis
- Music identification & analysis

3

Brute Force Exact Match

Pattern: **CCATT**

Text: **ACTGCCATTCCTTAGGGCCATGTG**

Algorithm?
Example of worst case?
Running time?

4

Rabin-Karp, 1981

Pattern: **CCATT**

Text: **ACTGCCATTCCTTAGGGCCATGTG**

Hash pattern & $|P|$ -substrings

- Compare hashes.
- Compare strings only when hashes match.
- $h(xb)$ easily computed from $h(ax)$, b

Best and worst cases?

5

Intuition for Better Algorithms

Pattern: **CCATT**

Text: **CCAGCCATTCCTTAGGGCCATGTG**

After failing match at 1st position, what should we do?

6

Quick Overview of Two Algorithms

Knuth-Morris-Pratt, 1977

- Use previous intuition.
- Preprocessing builds shift table based upon pattern prefixes. $O(|P|)$
- Each text character compared once or twice. $O(|T|)$

Boyer-Moore, 1977 & Turbo Boyer-Moore, 1992

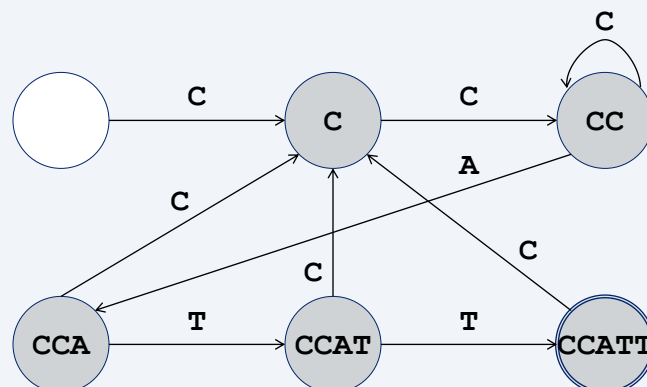
- Use previous intuition, along with similar heuristics. Complicated.
- Match from right end of pattern.
- Can often skip some text characters. $O(|T|)$, with $O\left(\frac{|T|}{|P|}\right)$ best case.

7

Finite Automata Matching Example

Pattern: **CCATT**

Text: **CCAGCCATTCCTTAGGGCCATGTG**



$O(|P| \cdot |\Sigma|)$ to build

8

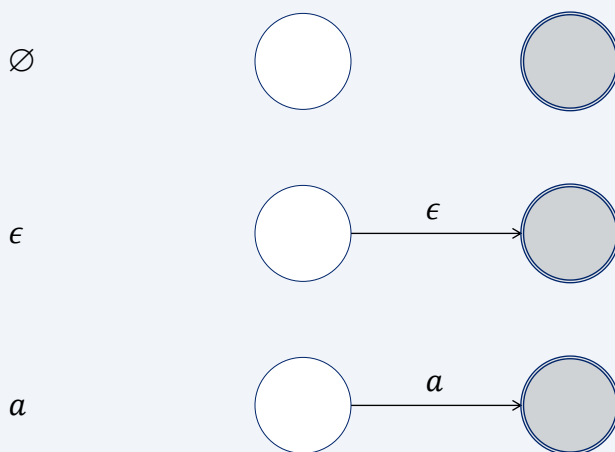
Regular Expressions

$(0 + 1)(00 + 01 + 10 + 11)^*$

Syntax: $\emptyset, \epsilon, a, rs, r + s, r^*$

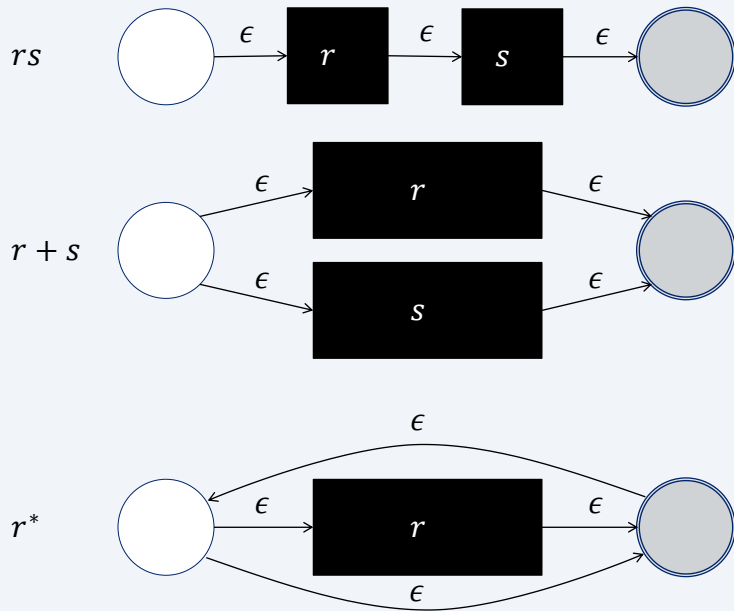
9

Regular Expression to Finite Automaton



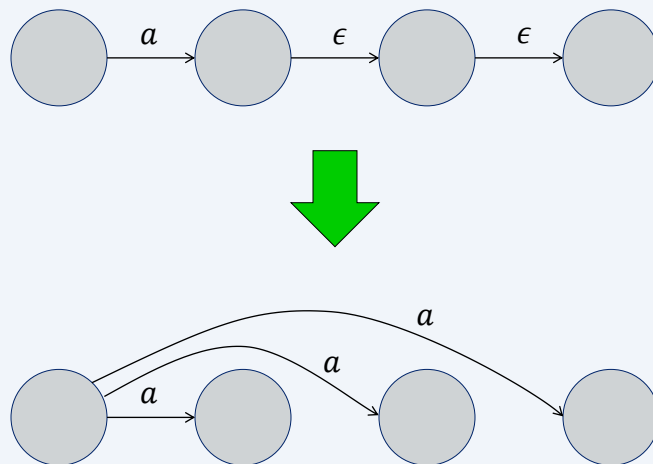
10

Regular Expression to Finite Automaton



11

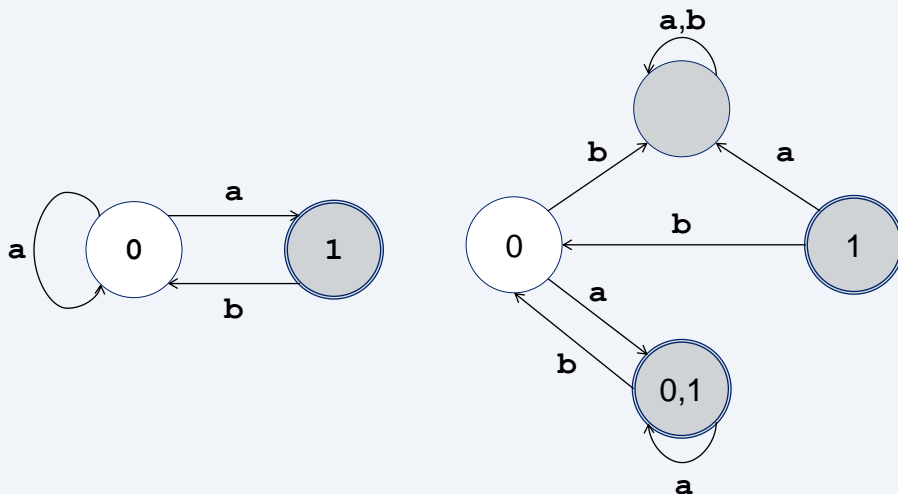
Eliminating ϵ



12

Eliminating Non-Determinism

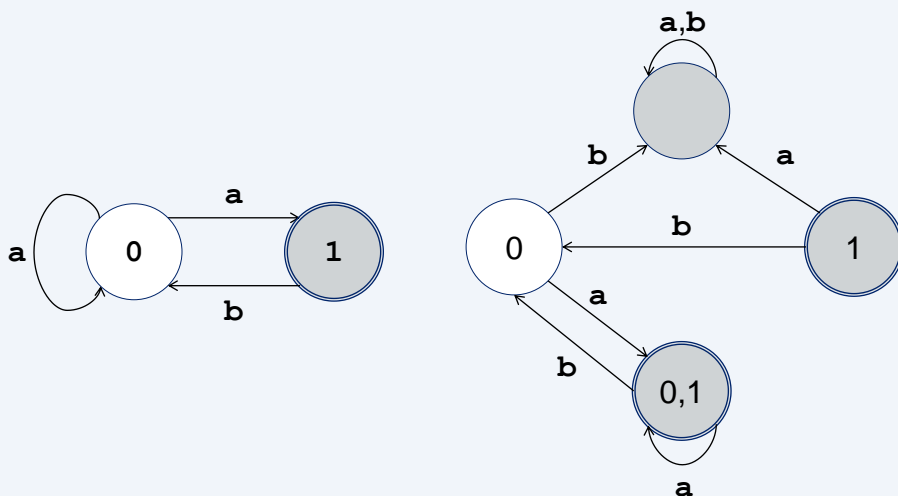
Each state is a set of the old states.



13

Can Minimize

States unreachable or equivalent? $O(|Q|^2)$



14

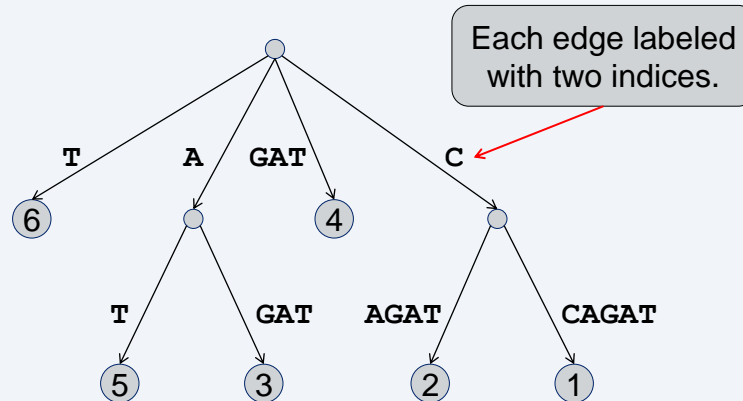
Finite Automaton Approach Summary

Preprocessing expensive for REs.
Matching is flexible and still linear.

15

Suffix Tree (Trie)

Text: **CCAGAT**

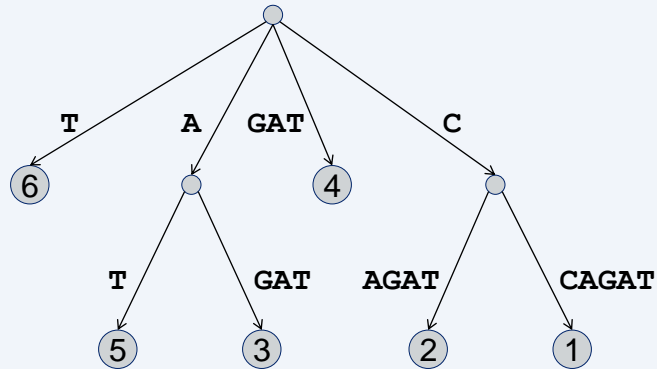


How many leaves, nodes, edges? Total space?
How to search for a pattern? Running time?

16

Require Suffixes to End at Leaves

Text: **CCAGAT**

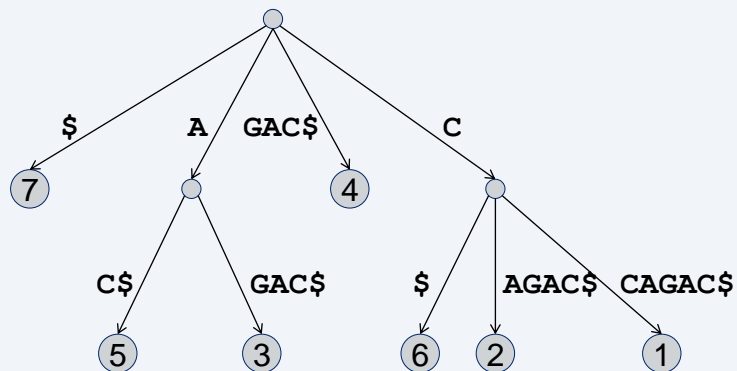


Reason: Simplicity. Distinguish nodes & leaves.
Example text that would break that?

17

Forcing Suffixes to End at Leaves

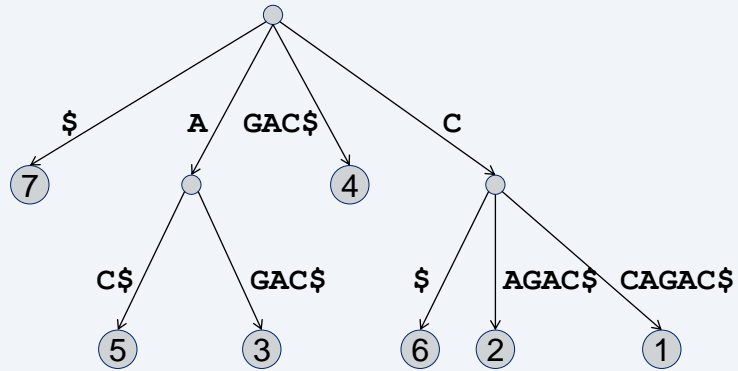
Text: **CCAGAC\$**



18

How Would You Create the Tree?

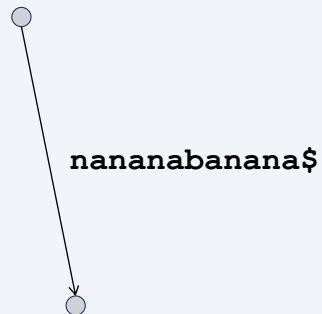
Text: CCAGAC\$



19

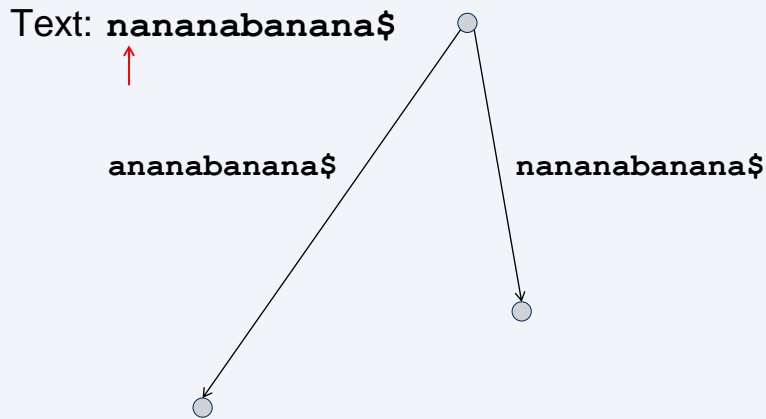
Suffix Tree Construction Example

Text: nananabanana\$



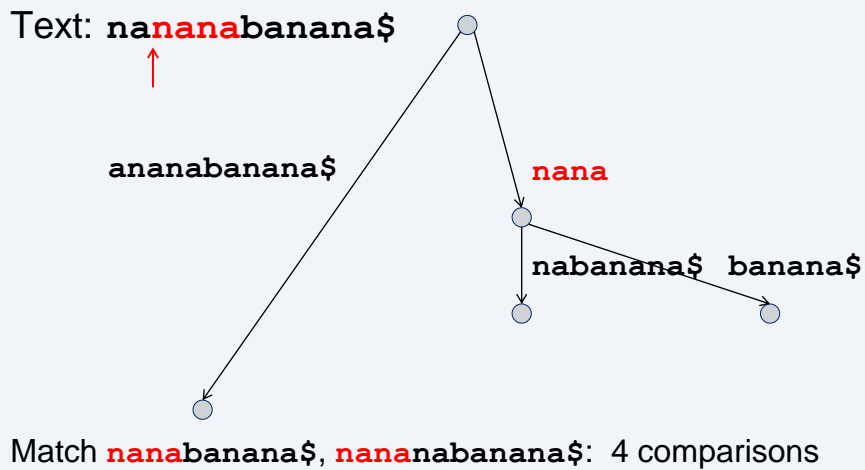
20

Suffix Tree Construction Example



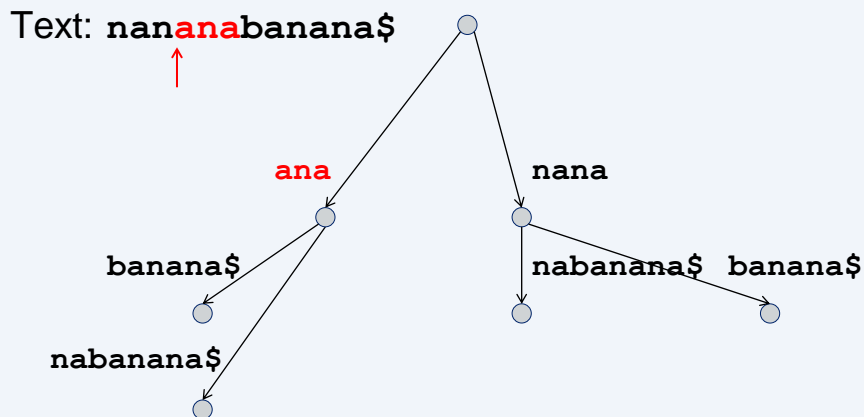
21

Suffix Tree Construction Example



22

Suffix Tree Construction Example



- Match **an**banana\$, **ana**nabananana\$: 3 redundant comps
- Next ...
- Match **na**banana\$, **na**nabananana\$: 2 redundant comps
- Match **a**banana\$, **a**nabananana\$: 1 redundant comp

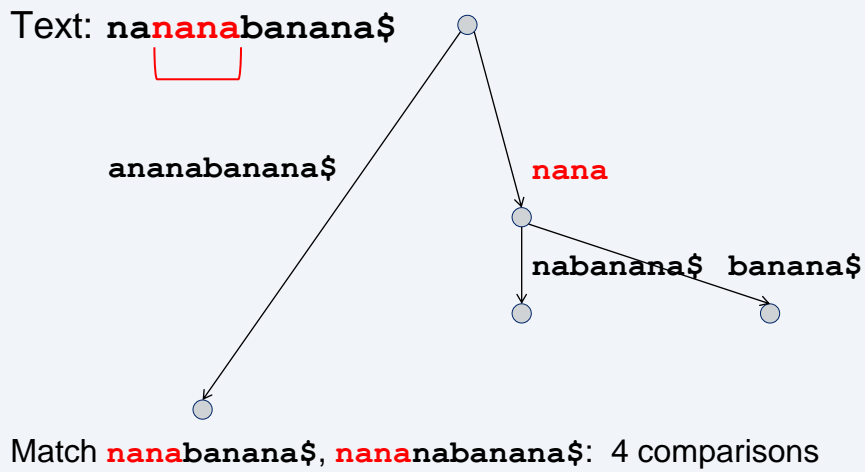
23

First Algorithm Improvement: No Redundant Matching

When inserting aXY , nanabananana\$
 If we discover aX is a prefix in tree, nana
 Then X will also be a prefix in tree. ana
 So, don't bother matching to verify.

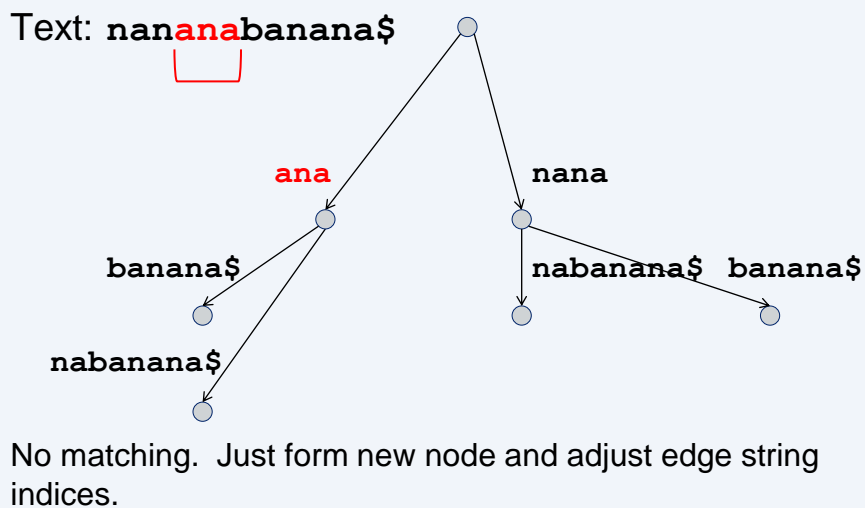
24

Suffix Tree Construction Example



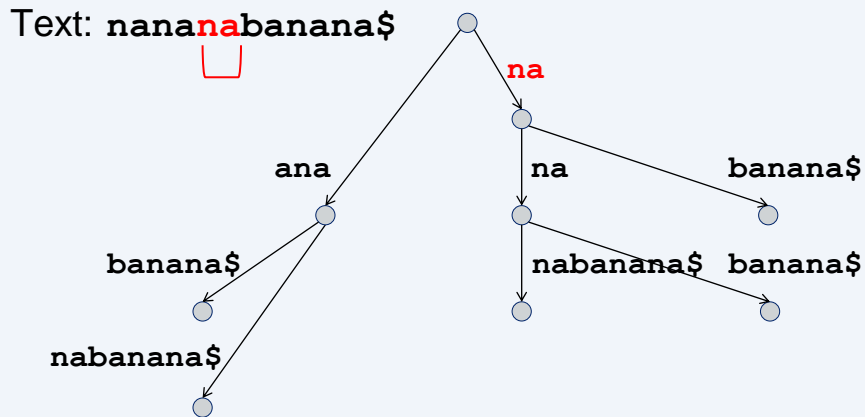
25

Suffix Tree Construction Example



26

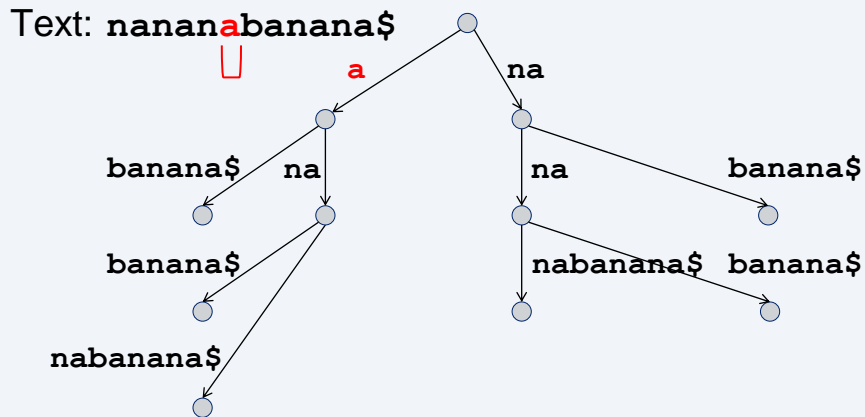
Suffix Tree Construction Example



No matching. Just form new node and adjust edge string indices.

27

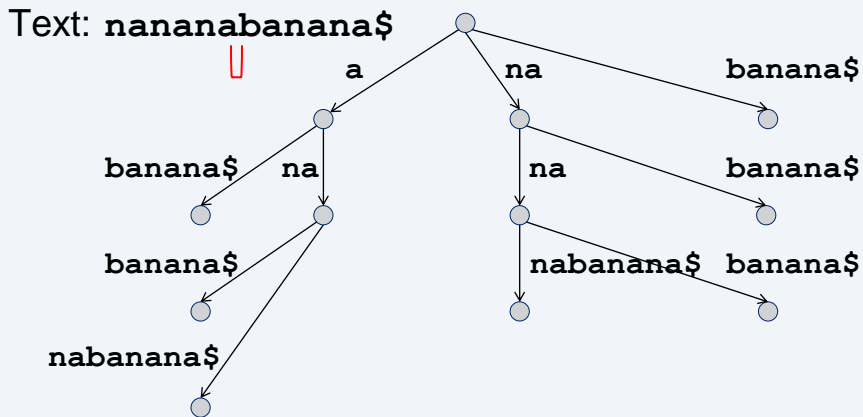
Suffix Tree Construction Example



No matching. Just form new node and adjust edge string indices.

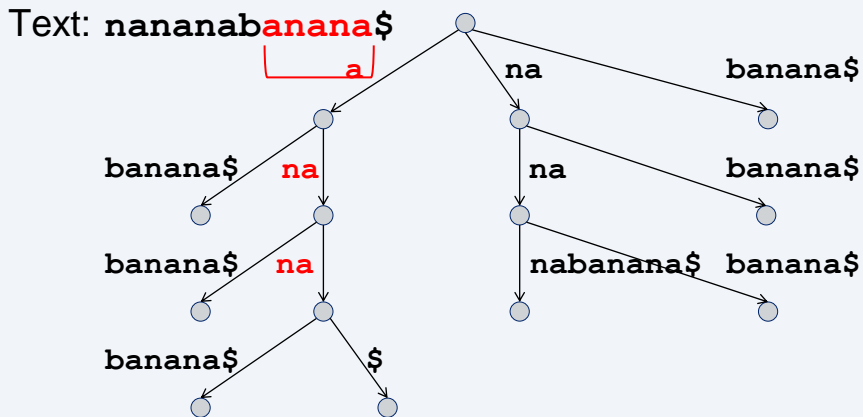
28

Suffix Tree Construction Example



29

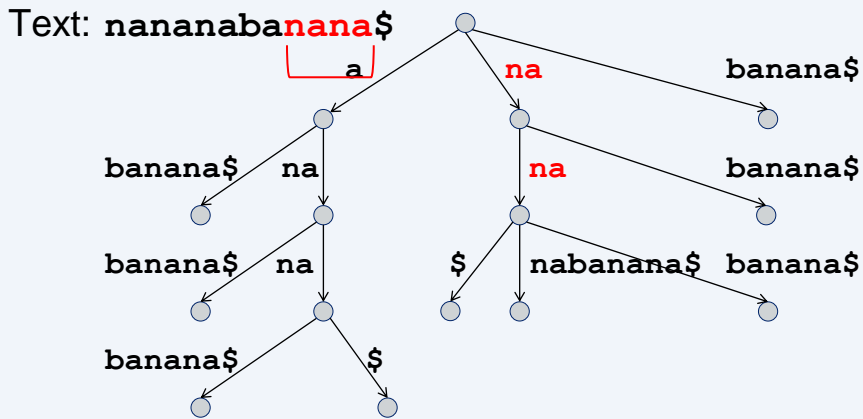
Suffix Tree Construction Example



anana is there, but it takes work to match & follow links.

30

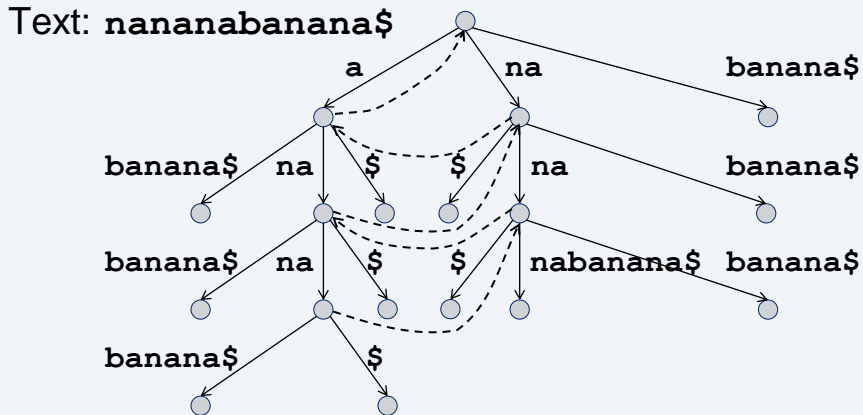
Suffix Tree Construction Example



nana is there, but it takes work to match & follow links.

31

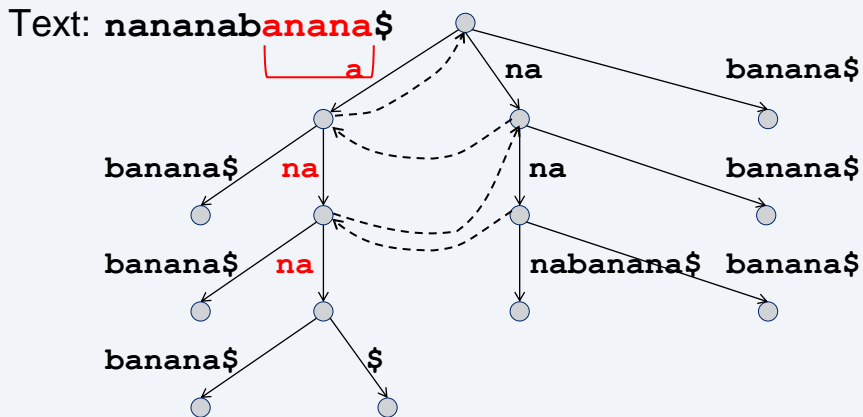
Second Algorithm Improvement: Suffix Links



Each internal node for xA has pointer to node for A .

32

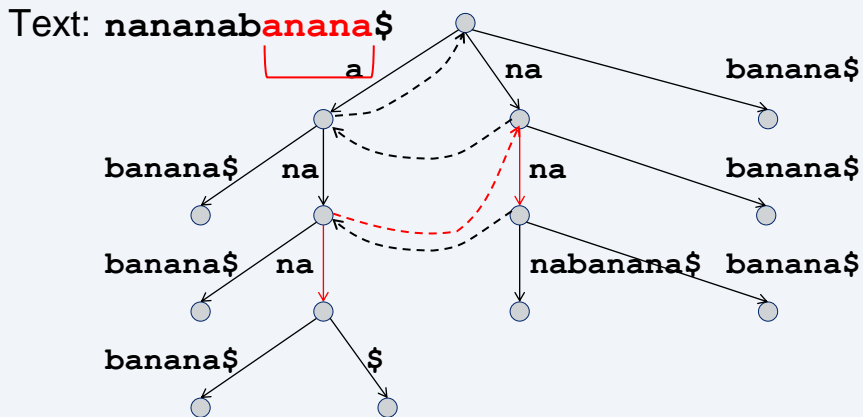
Suffix Tree Construction Example



Have to follow links for match on first time.
Need to create suffix link for new node.

33

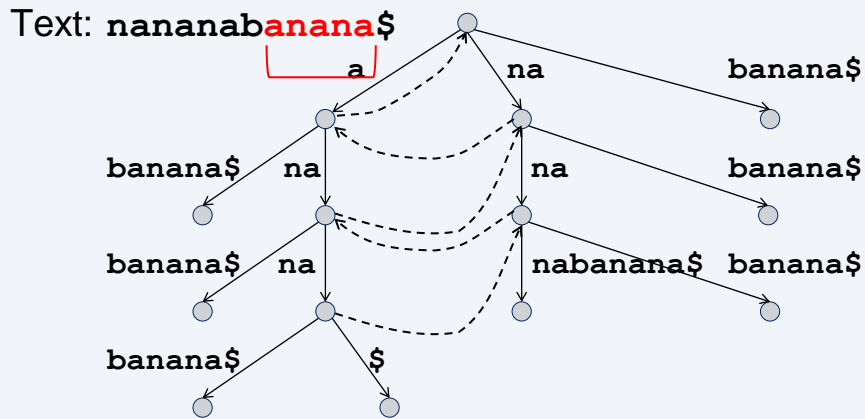
Suffix Tree Construction Example



Need to create suffix link for new node.
Use parent's suffix link to get close.

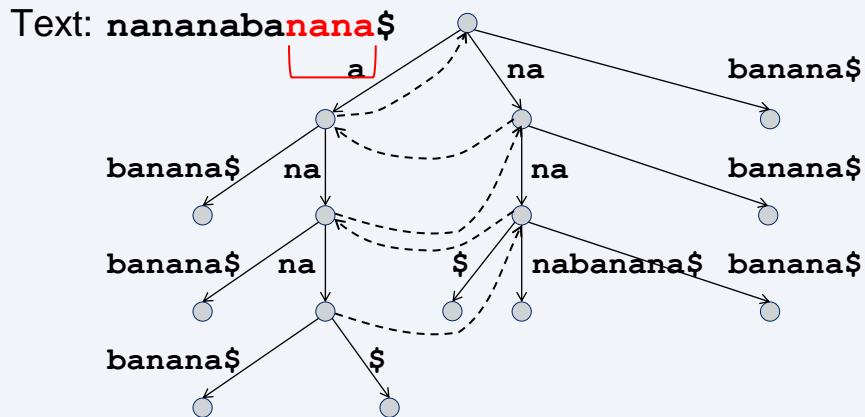
34

Suffix Tree Construction Example



35

Suffix Tree Construction Example



Already found node. No searching or matching.

36

Almost Correct Analysis of Construction

Two indices: $i \quad \underbrace{\quad}_j$

j : Each increment takes $O(1)$ time.

- Just search for one more character.

i : Each increment takes $O(1)$ time.

- Follow suffix link, or from root, get to suffix match.
- Possibly split node & create suffix link.
- Add one edge & leaf.

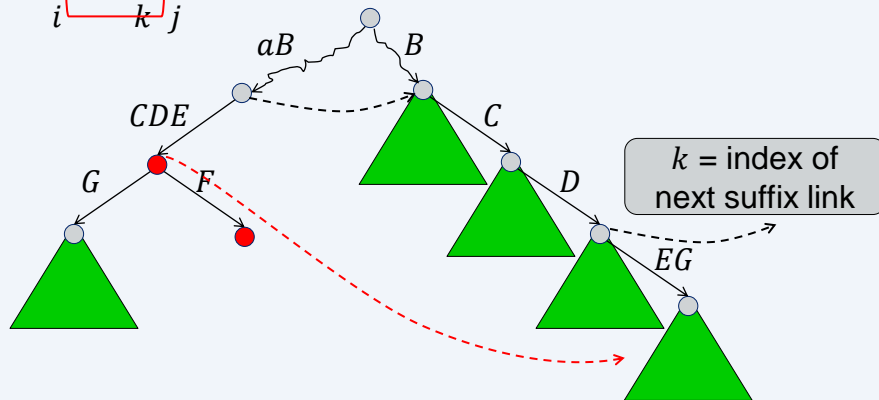
i, j each incremented n times $\rightarrow O(n)$ total.

41

Creating & Following Suffix Links

When creating new node, need new suffix link.
Follow parent's suffix link to get close, then search down.

Text: $ZaBCDEF$
 $i \quad \underbrace{\quad}_k \quad j$



42

Correct Analysis of Construction

Three indices: $i \text{---} k \text{---} j$ (k not part of alg.)

j : Each increment takes $O(1)$ time.

k : Follow some number l of links in $O(l)$ time.

Increments k by at least l .

i : Each increment takes $O(1)$ time in addition to that considered for k .

i, j, k each incremented at most n times $\rightarrow O(n)$ total.

43

Some Applications of Suffix Trees

Search for fixed patterns

Search for regular expressions

Find longest common substrings, longest repeated substrings

Find most commonly repeated substrings

Find maximal palindromes

Find Lempel-Ziv decomposition (for text compression)

As used in

Bioinformatics

Data compression

Data clustering

44

Supplementary Resources

Exact String Matching Algorithms

>30 algorithms, with animations

Course on string matching (Biosequencing)

Wikipedia: Suffix trees

Tutorial on Suffix trees

Tutorial on Suffix trees with applet & code

Notes on string matching and building suffix trees

Suffix tree slides adapted from those by Guy Blelloch, CMU 15-853.

45