



13 Text Indexing – Searching entire genomes

2 February 2026

Prof. Dr. Sebastian Wild

Learning Outcomes

Unit 13: *Text Indexing*

1. Know and understand methods for text indexing: *inverted indexes, suffix trees, (enhanced) suffix arrays*
2. Know and understand *generalized suffix trees*
3. Know properties, in particular *performance characteristics*, and limitations of the above data structures.
4. Design (simple) *algorithms based on suffix trees*.
5. Understand *construction algorithms* for suffix arrays.

Outline

13 Text Indexing

- 13.1 Inverted Indexes and Tries
- 13.2 Suffix Trees
- 13.3 Applications
- 13.4 Generalized Suffix Trees
- 13.5 Suffix Arrays
- 13.6 Linear-Time Suffix Sorting: Inducing Order
- 13.7 Linear-Time Suffix Sorting: The DC3 Algorithm

13.1 Inverted Indexes and Tries

Text indexing

- ▶ *Text indexing* (also: *offline text search*):
 - ▶ case of string matching: find $P[0..m]$ in $T[0..n]$
 - ▶ but with *fixed* text \rightsquigarrow preprocess T (instead of P)
 - \rightsquigarrow expect many queries P , answer them without looking at all of T
 - \rightsquigarrow essentially a data structuring problem: “building an *index* of T ”
- ▶ application areas
 - ▶ web search engines
 - ▶ online dictionaries
 - ▶ online encyclopedia
 - ▶ DNA/RNA data bases
 - ▶ . . . searching in any collection of text documents (that grows only moderately)

Latin: “one who points out”

Inverted indexes

≈ same as “indices”

- ▶ original indexes in books: list of (key) words ↪ page numbers where they occur
- ▶ assumption: searches are only for **whole** (key) **words**
- ~~ often reasonable for natural language text

Inverted indexes

≈ same as “indices”

- ▶ original indexes in books: list of (key) words ↪ page numbers where they occur
- ▶ assumption: searches are only for **whole** (key) words
- ~~ often reasonable for natural language text

Inverted index:

- ▶ collect all words in T
 - ▶ can be as simple as splitting T at whitespaces
 - ▶ actual implementations typically support *stemming* of words
 - goes → go, cats → cat
 - language specific!
- ▶ store mapping from words to a list of occurrences ~~ how?

Clicker Question



Do you know what a *trie* is?

- A** A what? No!
- B** I have heard the term, but don't quite remember.
- C** I remember hearing about it in a module.
- D** Sure.



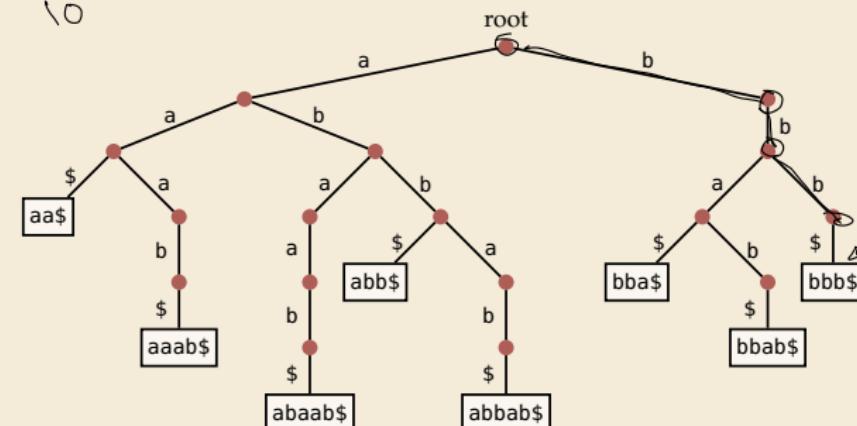
→ *sli.do/cs566*

Tries

- ▶ efficient dictionary data structure for strings
- ▶ name from retrieval, but pronounced “try”
- ▶ tree based on symbol comparisons
- ▶ **Assumption:** stored strings are *prefix-free* (no string is a prefix of another)
 - ▶ strings of same length ✓
 - ▶ some character $\notin \Sigma$
 - ▶ strings have “end-of-string” marker \$ ✓

Example:

{aa\$, aaab\$, abaab\$, abb\$,
abbab\$, bba\$, bbab\$, bbb\$}



Clicker Question

Suppose we have a trie that stores n strings over $\Sigma = \{A, \dots, Z\}$. Each stored string consists of m characters.

We now search for a query string Q with $|Q| = q$ (with $q \leq m$).
How many **nodes** in the trie are visited during this **query**?



A $\Theta(\log n)$

B $\Theta(\log(nm))$

C $\Theta(m \cdot \log n)$

D $\Theta(m + \log n)$

E $\Theta(m)$

F $\Theta(\log m)$

G $\Theta(q)$

H $\Theta(\log q)$

I $\Theta(q \cdot \log n)$

J $\Theta(q + \log n)$



→ *sli.do/cs566*

Clicker Question

Suppose we have a trie that stores n strings over $\Sigma = \{A, \dots, Z\}$. Each stored string consists of m characters.

We now search for a query string Q with $|Q| = q$ (with $q \leq m$).
How many **nodes** in the trie are visited during this **query**?



A ~~$\Theta(\log n)$~~

F ~~$\Theta(\log m)$~~

B ~~$\Theta(\log(mn))$~~

G $\Theta(q)$ ✓

C ~~$\Theta(m \log n)$~~

H ~~$\Theta(\log q)$~~

D ~~$\Theta(m + \log n)$~~

I ~~$\Theta(q \log n)$~~

E ~~$\Theta(m)$~~

J ~~$\Theta(q + \log n)$~~



→ *sli.do/cs566*

Clicker Question



Suppose we have a trie that stores n strings over $\Sigma = \{A, \dots, Z\}$. Each stored string consists of m characters.

How many **nodes** does the trie have **in total in the worst case?**

A $\Theta(n)$

B $\Theta(n + m)$

C $\Theta(n \cdot m)$

D $\Theta(n \log m)$

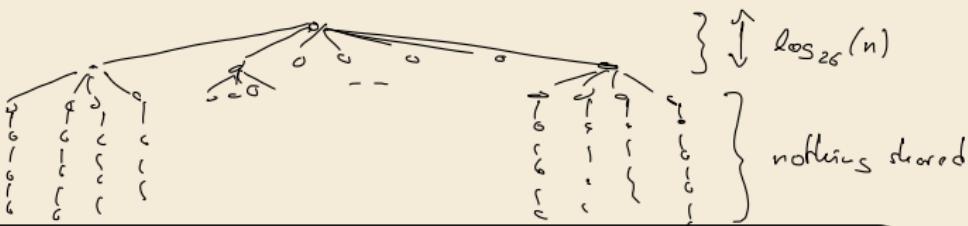
E $\Theta(m)$

F $\Theta(m \log n)$



→ *sli.do/cs566*

Clicker Question



Suppose we have a trie that stores n strings over $\Sigma = \{A, \dots, Z\}$. Each stored string consists of m characters.

How many **nodes** does the trie have **in total** in the worst case?



- A ~~$\Theta(n)$~~
- B ~~$\Theta(n + m)$~~ total # chars
- C $\Theta(n \cdot m)$ ✓
- D ~~$\Theta(n \log m)$~~
- E ~~$\Theta(m)$~~
- F ~~$\Theta(m \log n)$~~

$$n(m - \log_2 n) + n$$

node in trie needs array of child pointers

$\sim O(5)$ space per node

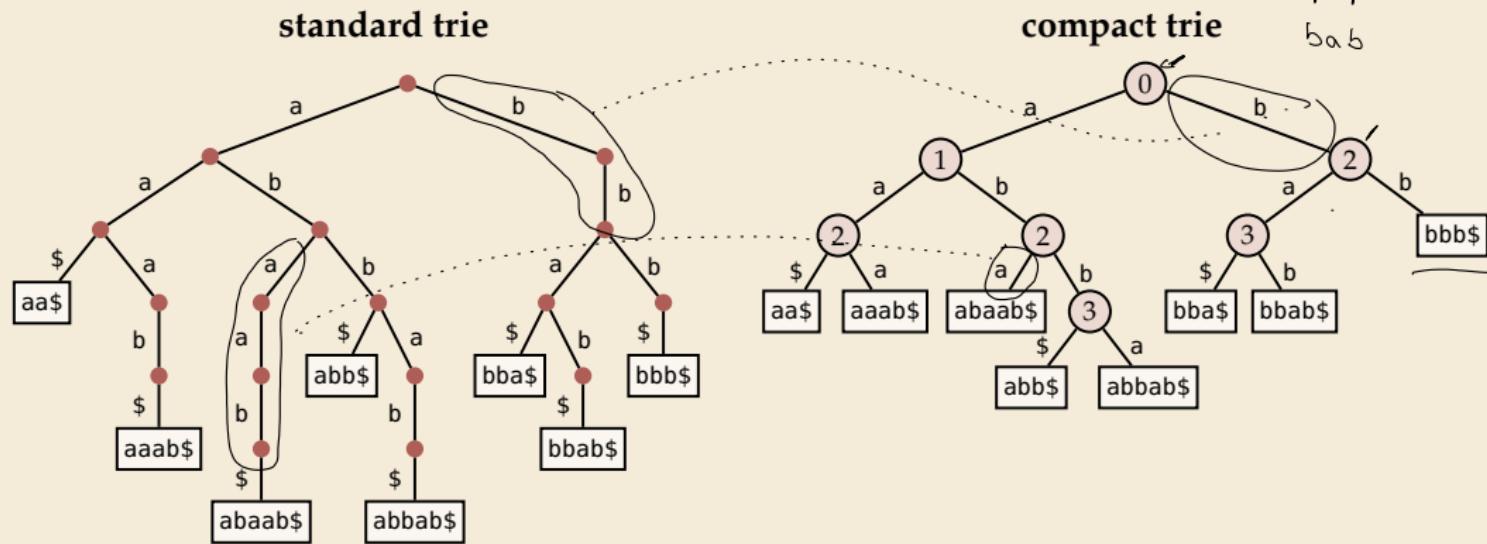


→ sli.do/cs566

Compact tries

- ▶ compress paths of unary nodes into single edge
- ▶ nodes store *index* of next character to check

=1 child



~~> searching slightly trickier, but same time complexity as in trie

- ▶ all nodes ≥ 2 children
- ~~> #nodes \leq #leaves = #strings
- ~~> linear space in $\# \text{strings}$

Tries as inverted index

 simple

 fast lookup

 cannot handle more general queries:

- ▶ search part of a word
- ▶ search phrase (sequence of words)

Tries as inverted index

👍 simple

👍 fast lookup

👎 cannot handle more general queries:

- ▶ search part of a word
- ▶ search phrase (sequence of words)

👎 what if the 'text' does not even have words to begin with?!

▶ biological sequences

```
ACAAGATGCCATTGCCCCGGCCCTCTGCTGCTGCTCTCGGGGCCACGGCACCGCTGCCCTGCCCTGGAGGGTGGCCCCACCGGC  
CGAGACAGCGAGCATATGCAGGAAGCGGCAGGAATAAGGAAAAGCAGCCTCTGACTTCTCGTTGGTGGTTGAGTGGACCTCCAGGC  
CAGTGCCGGCCCCCTCATAGGAGAGGAAGCTGGGAGGTGGCAGGCGCAGGAAGGCGCACCCCCCAGCAATCCGCGCCGGGACAGAA  
TGCCCTGCAGGAACTTCTTCTGGAAGACCTCTCCTGCAAATAAACCTCACCATGAATGCTCACGCAAGTTAATTACAGACCTGAA
```

▶ binary streams

```
0000001010100111010111000001111100011111011111001101101000011100010011011110000010001101010  
0110110000110101101000000010000000011101011000001000011110101110110010001100101101110111111  
110001010001011001010000001110101010011000000001101100001100111110000101 0101011101111000011  
101011100100101010100000111110100110000001111001101010000000100100100000101100011000110111
```

↝ need new ideas

13.2 Suffix Trees

Suffix trees – A ‘magic’ data structure

Appetizer: Longest common substring problem

- ▶ Given: strings S_1, \dots, S_k **Example:** $S_1 = \text{superiorcalifornialives}$, $S_2 = \underline{\text{sealiver}}$
- ▶ Goal: find the longest substring that occurs in all k strings

Suffix trees – A ‘magic’ data structure

Appetizer: Longest common substring problem

- Given: strings S_1, \dots, S_k Example: $S_1 = \text{superiorcalifornialives}, S_2 = \text{sealiver}$
 - Goal: find the longest substring that occurs in all k strings \rightsquigarrow alive



Can we do this in time $O(|S_1| + \dots + |S_k|)$? How??

Suffix trees – A ‘magic’ data structure

Appetizer: Longest common substring problem

- ▶ Given: strings S_1, \dots, S_k
- Example: $S_1 = \text{superiorcalifornialives}$, $S_2 = \text{sealiver}$
- ▶ Goal: find the longest substring that occurs in all k strings
- ~~~ alive



Can we do this in time $O(|S_1| + \dots + |S_k|)$? How??

Enter: *suffix trees*

- ▶ versatile data structure for index with full-text search
- ▶ linear time (for construction) and linear space
- ▶ allows efficient solutions for many advanced string problems



“Although the longest common substring problem looks trivial now, given our knowledge of suffix trees, it is very interesting to note that in 1970 Don Knuth conjectured that a linear-time algorithm for this problem would be impossible.”

[Gusfield: *Algorithms on Strings, Trees, and Sequences* (1997)]

Suffix trees – Definition

- suffix tree \mathcal{T} for text $T = T[0..n)$ = compact trie of all suffixes of $T\$$ (set $T[n] := \$$)

Suffix trees – Definition

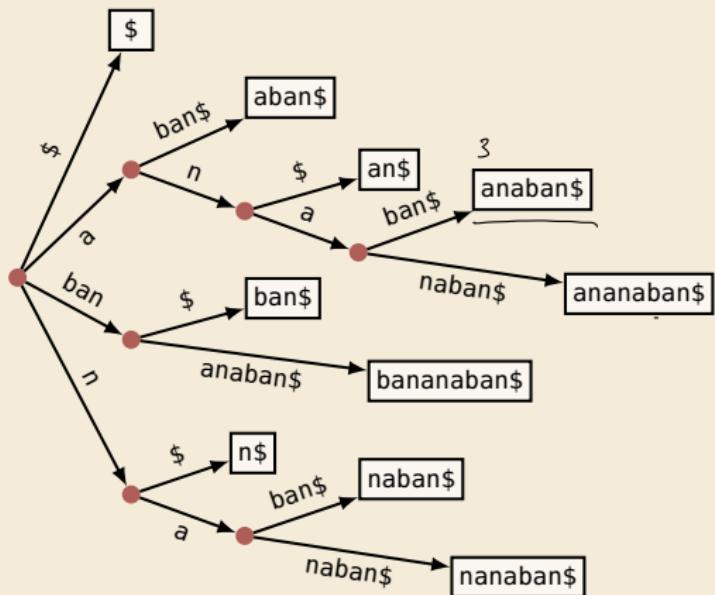
- suffix tree \mathcal{T} for text $T = T[0..n)$ = compact trie of all suffixes of $T\$$ (set $T[n] := \$$)

Example:

$T = \text{bananaban\$}$

suffixes: $\{\text{bananaban\$}, \text{ananaban\$}, \text{nanaban\$}, \text{anaban\$}, \text{naban\$}, \text{aban\$}, \text{ban\$}, \text{an\$}, \text{n\$}, \$\}$

$T = \boxed{\begin{matrix} 0 & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 \\ \text{b} & \text{a} & \text{n} & \text{a} & \text{n} & \text{a} & \text{b} & \text{a} & \text{n} & \$ \end{matrix}}$



Suffix trees – Definition

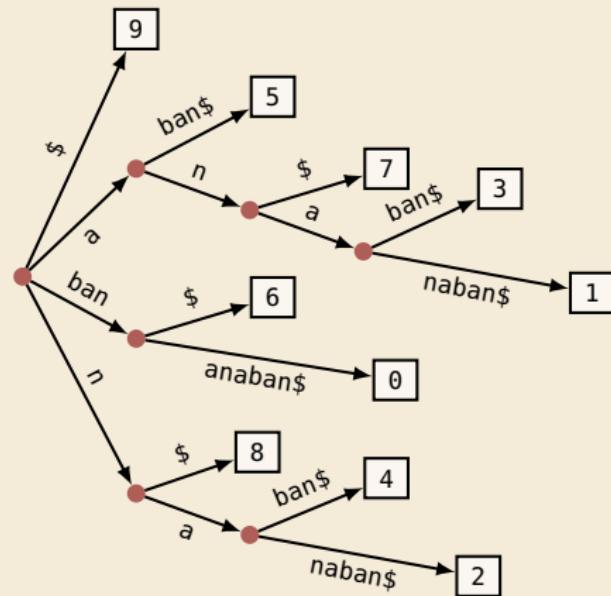
- suffix tree \mathcal{T} for text $T = T[0..n)$ = compact trie of all suffixes of $T\$$ (set $T[n] := \$$)
- except: in leaves, store *start index* (instead of copy of actual string)

Example:

$T = \text{bananaban\$}$

suffixes: { $\text{bananaban\$}$, $\text{ananaban\$}$, $\text{nanaban\$}$,
 $\text{anaban\$}$, $\text{nab\$}$, $\text{aban\$}$, $\text{ban\$}$, $\text{an\$}$, $\text{n\$}$, $\$\$$ }

$T = \begin{array}{cccccccccc} 0 & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 \\ \boxed{\mathbf{b}} & \boxed{\mathbf{a}} & \boxed{\mathbf{n}} & \boxed{\mathbf{a}} & \boxed{\mathbf{n}} & \boxed{\mathbf{a}} & \boxed{\mathbf{b}} & \boxed{\mathbf{a}} & \boxed{\mathbf{n}} & \boxed{\$} \end{array}$



Suffix trees – Definition

- suffix tree \mathcal{T} for text $T = T[0..n)$ = compact trie of all suffixes of $T\$$ (set $T[n] := \$$)
- except: in leaves, store *start index* (instead of copy of actual string)

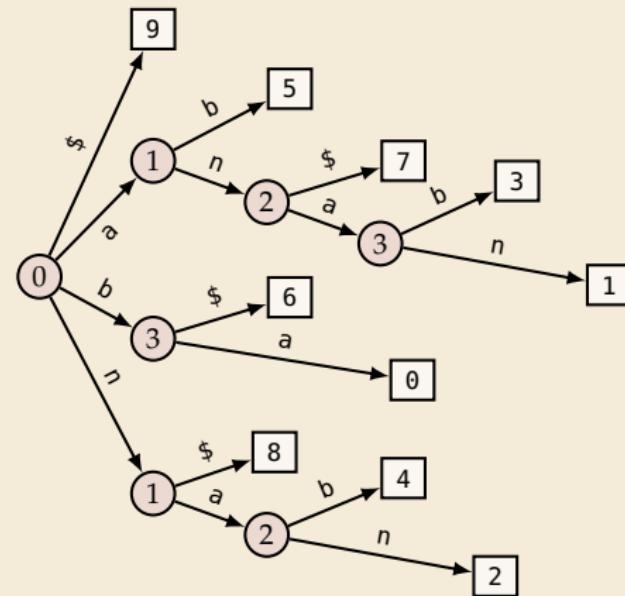
Example:

$T = \text{bananaban\$}$

suffixes: {bananaban\$, ananaban\$, nanaban\$,
ananabn\$, naban\$, aban\$, ban\$, an\$, n\$, \$}

$T = \begin{array}{cccccccccc} 0 & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 \\ | & | & | & | & | & | & | & | & | & | \\ \boxed{\text{b}} & \boxed{\text{a}} & \boxed{\text{n}} & \boxed{\text{a}} & \boxed{\text{n}} & \boxed{\text{a}} & \boxed{\text{b}} & \boxed{\text{a}} & \boxed{\text{n}} & \boxed{\$} \end{array}$

- also: edge labels like in compact trie
- (more readable form on slides to explain algorithms)



Suffix trees – Construction

- ▶ $T[0..n]$ has $n + 1$ suffixes (starting at character $i \in [0..n]$)
- ▶ We can build the suffix tree by inserting each suffix of T into a compressed trie.
But that takes time $\Theta(n^2)$. ↵ not interesting!

Suffix trees – Construction

- ▶ $T[0..n]$ has $n + 1$ suffixes (starting at character $i \in [0..n]$)
- ▶ We can build the suffix tree by inserting each suffix of T into a compressed trie. But that takes time $\Theta(n^2)$. ↵ not interesting!



same order of growth as reading the text!
Amazing result: Can construct the suffix tree of T in $\Theta(n)$ time!

- ▶ algorithms are a bit tricky to understand → EAA Buch
- ▶ but were a theoretical breakthrough
- ▶ and they are efficient in practice (and heavily used)!

↪ for now, take linear-time construction for granted. What can we do with them?

Clicker Question



Recap: Check all correct statements about suffix tree \mathcal{T} of $T[0..n]$.

- A** We require T to end with $\$$.
- B** The size of \mathcal{T} can be $\Omega(n^2)$ in the worst case.
- C** \mathcal{T} is a standard trie of all suffixes of $T\$$.
- D** \mathcal{T} is a compact trie of all suffixes of $T\$$.
- E** The leaves of \mathcal{T} store (a copy of) a suffix of $T\$$.
- F** Naive construction of \mathcal{T} takes $\Omega(n^2)$ (worst case).
- G** \mathcal{T} can be computed in $O(n)$ time (worst case).
- H** \mathcal{T} has n leaves.



→ *sli.do/cs566*

Clicker Question



Recap: Check all correct statements about suffix tree \mathcal{T} of $T[0..n]$.

- A** We require T to end with \$. ✓
- B** ~~The size of \mathcal{T} can be $\Omega(n^2)$ in the worst case.~~
- C** ~~\mathcal{T} is a standard trie of all suffixes of $T\$$.~~
- D** \mathcal{T} is a compact trie of all suffixes of $T\$$. ✓
- E** ~~The leaves of \mathcal{T} store (a copy of) a suffix of $T\$$.~~
- F** Naive construction of \mathcal{T} takes $\Omega(n^2)$ (worst case). ✓
- G** \mathcal{T} can be computed in $O(n)$ time (worst case). ✓
- H** ~~\mathcal{T} has n leaves. $n \rightarrow 1$~~



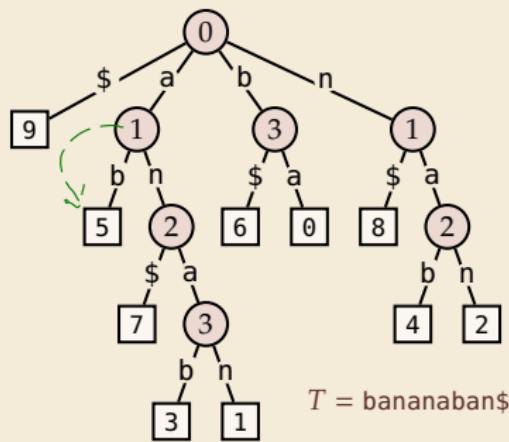
→ sli.do/cs566

13.3 Applications

Applications of suffix trees

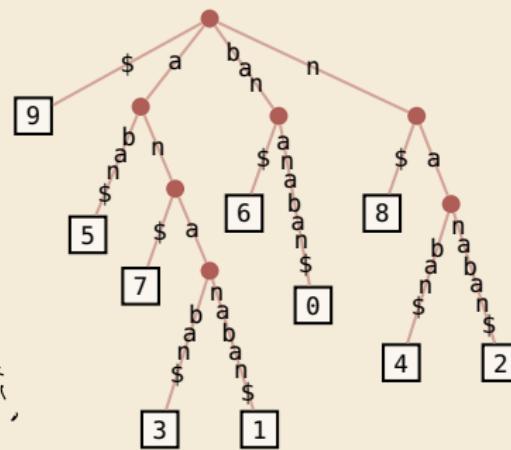
- In this section, always assume suffix tree \mathcal{T} for T given.

Recall: \mathcal{T} stored like this:



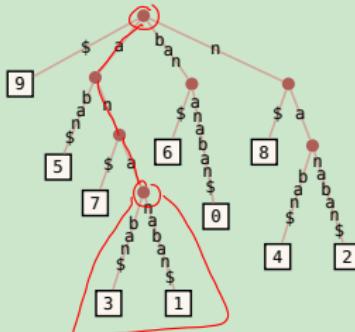
$T = \text{bananaban\$}$

but think about this:



- Moreover: assume internal nodes store pointer to leftmost leaf in subtree.
- Notation: $T_i = T[i..n]$ (including \$)

Clicker Question



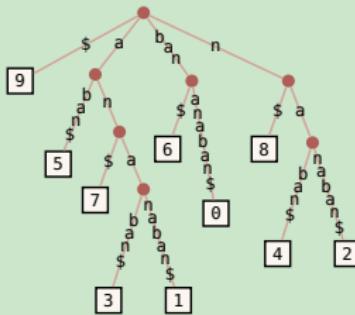
What does T 's suffix tree (on the left) tell you about the question whether T contains the pattern $P = \underline{ana}$?
Check all that apply to this example.

- A** Nothing.
- B** P occurs in T . ↪
- C** P does not occur in T .
- D** P occurs once in T .
- E** P occurs twice in T . ↫
- F** P starts at index 0.
- G** P starts at index 1.
- H** P starts at index 2.
- I** P starts at index 3.
- J** P starts at index 4.
- K** P starts at index 7.



→ *sli.do/cs566*

Clicker Question



What does T 's suffix tree (on the left) tell you about the question whether T contains the pattern $P = \text{ana}$?

Check all that apply to this example.

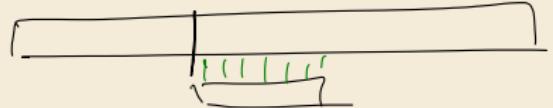
- A** ~~P starts at index 0.~~ **Nothing.**
- B** P occurs in T . ✓
- C** ~~P does not occur in T.~~ **P does not occur in T .**
- D** ~~P occurs once in T.~~ **P occurs once in T .**
- E** P occurs twice in T . ✓
- F** ~~P starts at index 0.~~ **P starts at index 0.**
- G** P starts at index 1. ✓
- H** ~~P starts at index 2.~~ **P starts at index 2.**
- I** P starts at index 3. ✓
- J** ~~P starts at index 4.~~ **P starts at index 4.**
- K** ~~P starts at index 7.~~ **P starts at index 7.**



→ *sli.do/cs566*

Application 1: Text Indexing / String Matching

- P occurs in $T \iff P$ is a prefix of a suffix of T
- we have all suffixes in \mathcal{T} !



Application 1: Text Indexing / String Matching

- P occurs in $T \iff P$ is a prefix of a suffix of T

- we have all suffixes in \mathcal{T} !

↷ (try to) follow path with label P , until

- ## 1. we get stuck

at internal node (no node with next character of P)

or *inside edge* (mismatch of next characters)

$\rightsquigarrow P$ does not occur in T

- ## **2. we run out of pattern**

reach end of P at internal node v or inside edge towards v

$\rightsquigarrow P$ occurs at all leaves in subtree of v

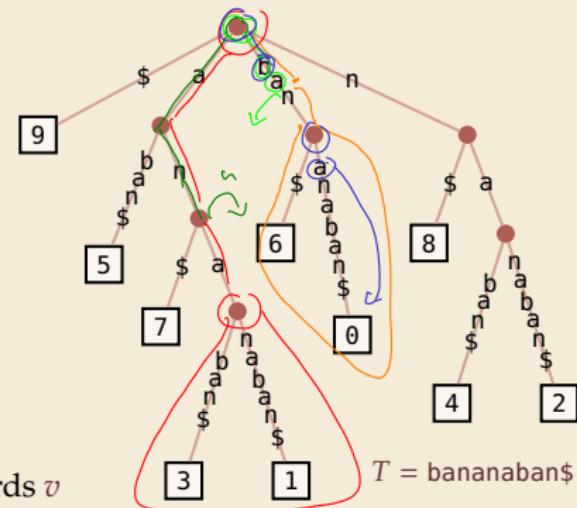
3. we run out of tree

reach a leaf ℓ with part of P left \rightsquigarrow compare P to ℓ .



This cannot happen when testing edge labels since $\$ \notin \Sigma$,
but needs check(s) in compact trie implementation!

- ▶ Finding first match (or NO_MATCH) takes $O(|P|)$ time!



Examples:

- ▶ $P = \underline{\text{ann}}$
 - ▶ $P = \underline{\text{baa}}$
 - ▶ $P = \underline{\text{ana}}$
 - ▶ $P = \underline{\text{ba}}$
 - ▶ $P = \underline{\text{brian}}$

Application 2: Longest repeated substring

► **Goal:** Find longest substring $T[i..i + \ell]$ that occurs also at $j \neq i$: $T[j..j + \ell] = T[i..i + \ell]$.



e.g. for compression \rightsquigarrow Unit 7

How can we efficiently check *all possible substrings*?

Application 2: Longest repeated substring

- Goal: Find longest substring $T[i..i + \ell]$ that occurs also at $j \neq i$: $T[j..j + \ell] = T[i..i + \ell]$.



e.g. for compression \rightsquigarrow Unit 7



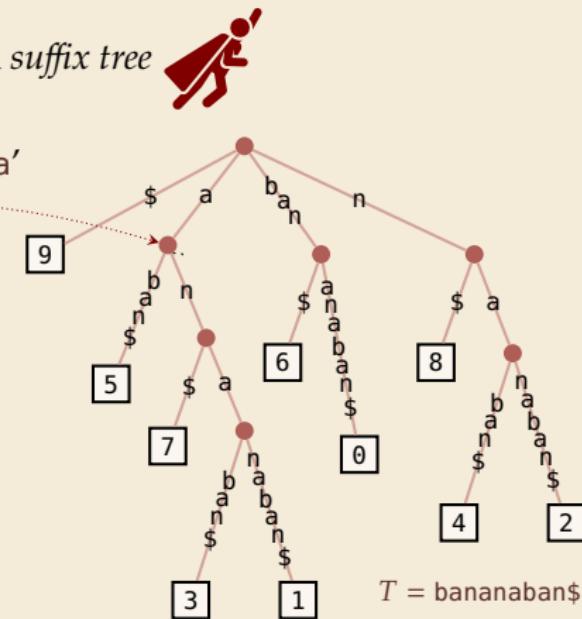
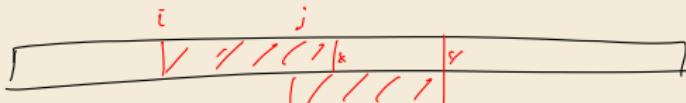
Repeated substrings = shared paths in *suffix tree*



- $T_5 = \text{aban\$}$ and $T_7 = \text{an\$}$ have *longest common prefix* 'a'

$\rightsquigarrow \exists$ internal node with path label 'a'

here single edge, can be longer path



$T = \text{bananaban\$}$

Application 2: Longest repeated substring

- Goal: Find longest substring $T[i..i + \ell]$ that occurs also at $j \neq i$: $T[j..j + \ell] = T[i..i + \ell]$.



e.g. for compression \rightsquigarrow Unit 7



Repeated substrings = shared paths in *suffix tree*



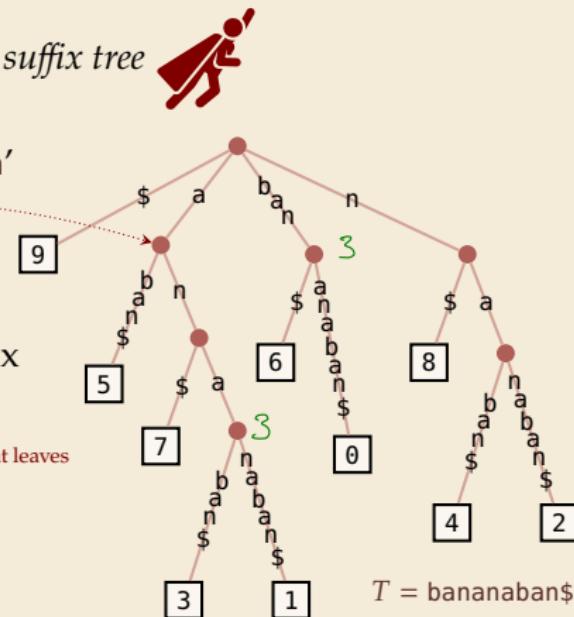
- $T_5 = \text{aban\$}$ and $T_7 = \text{an\$}$ have *longest common prefix* 'a'

$\rightsquigarrow \exists$ internal node with path label 'a'

here single edge, can be longer path

\rightsquigarrow longest repeated substring = longest common prefix
(LCP) of two suffixes

actually: adjacent leaves



$T = \text{bananaban\$}$

Application 2: Longest repeated substring

- Goal: Find longest substring $T[i..i + \ell]$ that occurs also at $j \neq i$: $T[j..j + \ell] = T[i..i + \ell]$.



e.g. for compression \rightsquigarrow Unit 7



Repeated substrings = shared paths in *suffix tree*



- $T_5 = \text{aban\$}$ and $T_7 = \text{an\$}$ have *longest common prefix* 'a'

$\rightsquigarrow \exists$ internal node with path label 'a'

here single edge, can be longer path

\rightsquigarrow longest repeated substring = longest common prefix
(LCP) of two suffixes

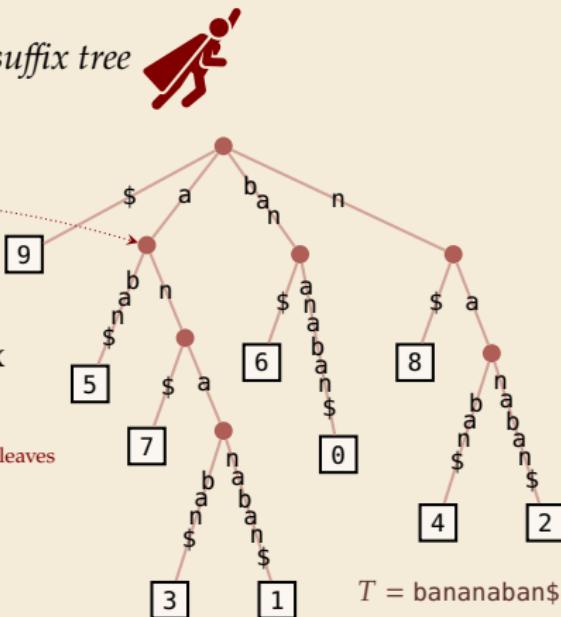
- Algorithm:

1. Compute string depth (=length of path label) of nodes

2. Find internal nodes with maximal string depth

actually: adjacent leaves

$\rightsquigarrow \Theta(n)$ time



13.4 Generalized Suffix Trees

Generalized suffix trees

- ▶ longest *repeated* substring (of one string) feels very similar to
longest *common* substring of several strings $T^{(1)}, \dots, T^{(k)}$ with $T^{(j)} \in \Sigma^{n_j}$
- ▶ can we solve that in the same way?
- ▶ could build the suffix tree for each $T^{(j)}$... but doesn't seem to help

Generalized suffix trees

- ▶ longest *repeated* substring (of one string) feels very similar to
longest *common* substring of several strings $T^{(1)}, \dots, T^{(k)}$ with $T^{(j)} \in \Sigma^{n_j}$
- ▶ can we solve that in the same way?
- ▶ could build the suffix tree for each $T^{(j)}$... but doesn't seem to help
- ~~ need a *single/joint* suffix tree for *several* texts

Generalized suffix trees

- ▶ longest *repeated* substring (of one string) feels very similar to
longest *common* substring of several strings $T^{(1)}, \dots, T^{(k)}$ with $T^{(j)} \in \Sigma^{n_j}$
- ▶ can we solve that in the same way?
- ▶ could build the suffix tree for each $T^{(j)}$... but doesn't seem to help
 - ~~ need a *single/joint* suffix tree for *several* texts

Enter: *generalized suffix tree*

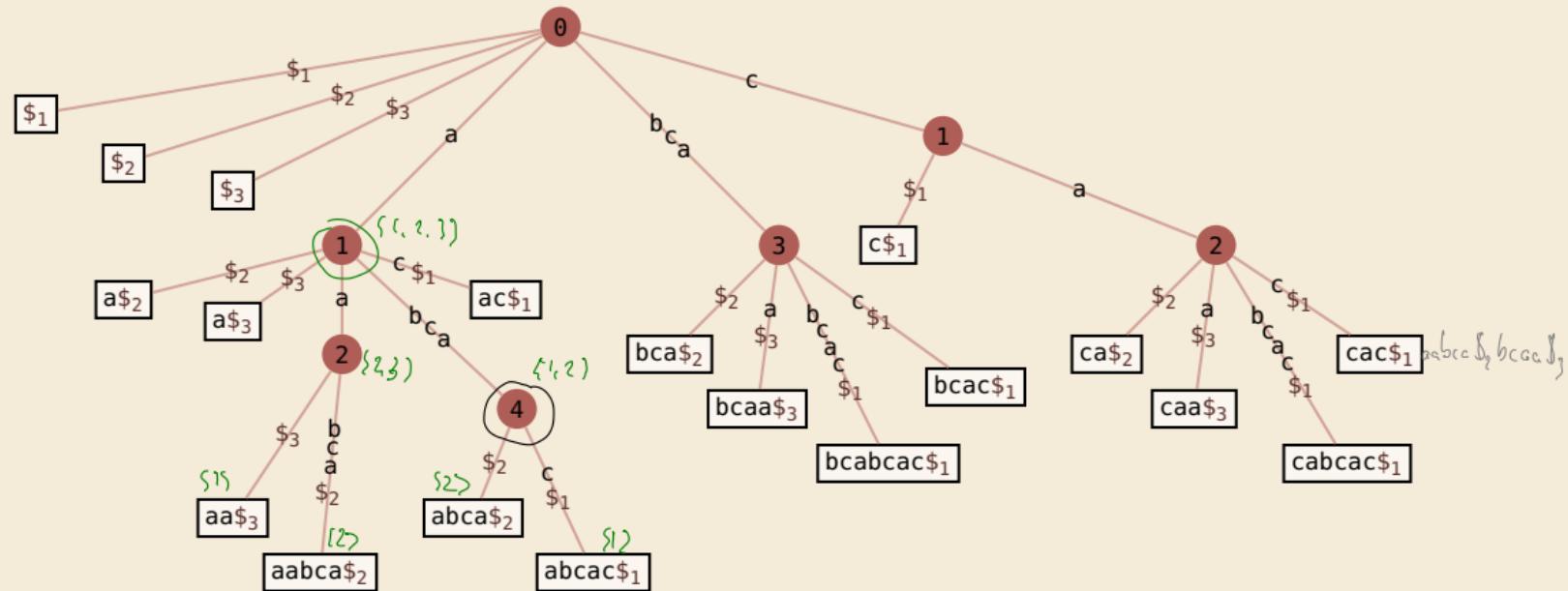
- ▶ Define $T := \underbrace{T^{(1)}\$_1 T^{(2)}\$_2 \dots T^{(k)}\$_k}_{\text{for } k \text{ new end-of-word symbols}}$
- ▶ Construct suffix tree \mathcal{T} for T
 - ~~ $\$_j$ -edges always leads to leaves ~~~ \exists leaf (j, i) for each suffix $T_i^{(j)} = T^{(j)}[i..n_j]$



Generalized Suffix Tree – Example

$$T^{(1)} = \text{bcabca}c, \quad T^{(2)} = \text{aabca}, \quad T^{(3)} = \text{bcaa}$$

$$T = \text{bcabca}c\$, \text{aabca}\$, \text{bcaa}\$$$



Application 3: Longest common substring

- With that new idea, we can find longest common substrings:

1. Compute generalized suffix tree \mathcal{T} .
2. Store with each node the *subset of strings* that contain its path label:
 - 2.1. Traverse \mathcal{T} bottom-up.
 - 2.2. For a leaf (j, i) , the subset is $\{j\}$.
 - 2.3. For an internal node, the subset is the union of its children.
3. In top-down traversal, compute *string depths* of nodes. (as above)
4. Report deepest node (by string depth) whose subset is $\underline{\{1, \dots, k\}}$.

- Each step takes time $\Theta(\underline{n})$ for $\underline{n} = n_1 + \dots + n_k$ the total length of all texts.

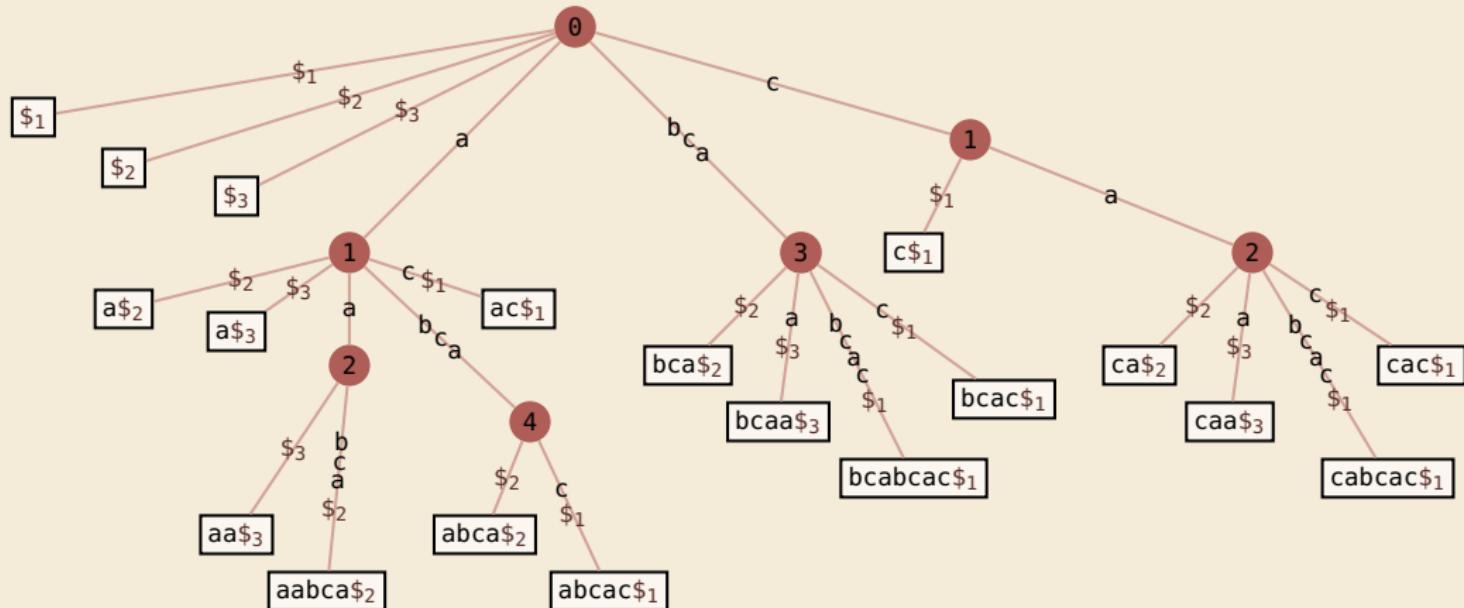
(for constant k)

“Although the longest common substring problem looks trivial now, given our knowledge of suffix trees, it is very interesting to note that in 1970 Don Knuth conjectured that a linear-time algorithm for this problem would be impossible.”

[Gusfield: Algorithms on Strings, Trees, and Sequences (1997)]

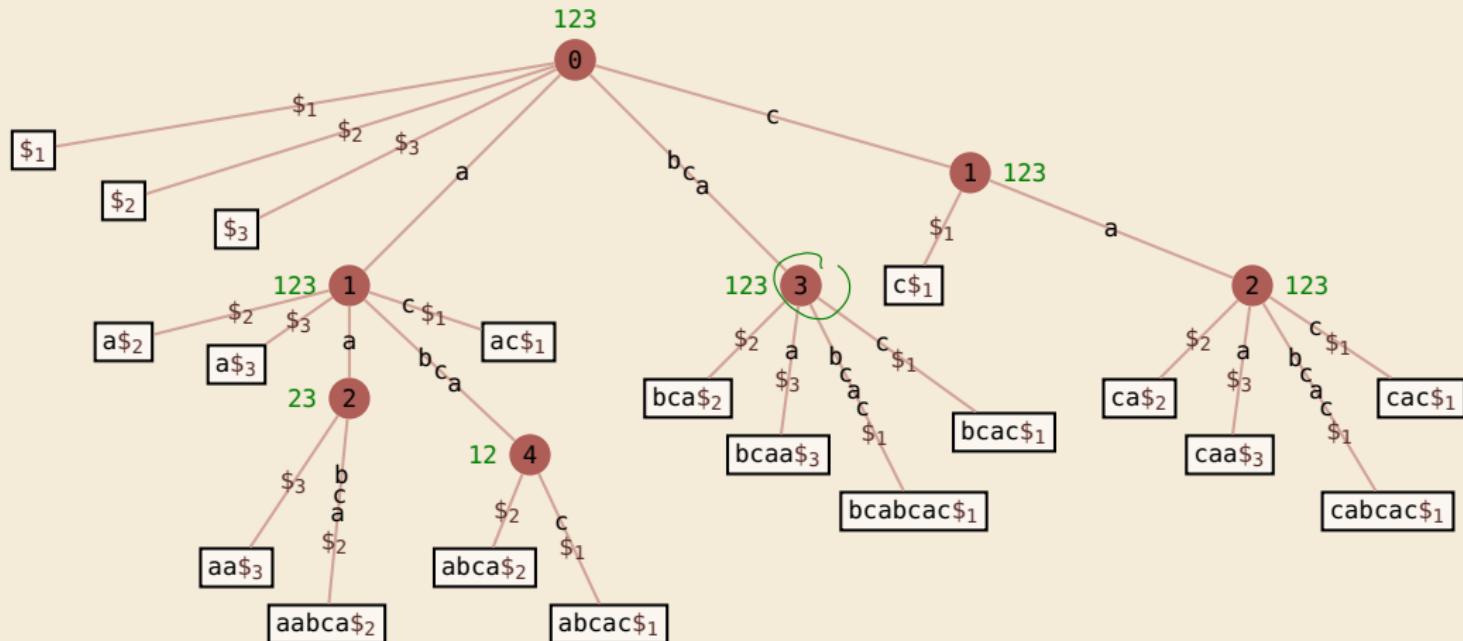
Longest common substring – Example

$$T^{(1)} = \text{bcabca}c, \quad T^{(2)} = \text{aabca}, \quad T^{(3)} = \text{bcaa}$$



Longest common substring – Example

$$T^{(1)} = \underline{bcab}c\underline{cac}, \quad T^{(2)} = \underline{aab}bc\underline{a}, \quad T^{(3)} = \underline{bcaa}$$



Further applications

Many more applications, in particular for problems on biological sequences

20+ described in Gusfield, *Algorithms on strings, trees, and sequences* (1999)

- ▶ key ingredient: longest common extensions

If you want to see more, come to *Algorithms of Bioinformatics* 😊

Suffix trees – Discussion

- ▶ Suffix trees were a threshold invention
 - ◀ linear time and space
 - ◀ suddenly many questions efficiently solvable in theory



Suffix trees – Discussion

- ▶ Suffix trees were a threshold invention
 - ◀ linear time and space
 - ◀ suddenly many questions efficiently solvable in theory



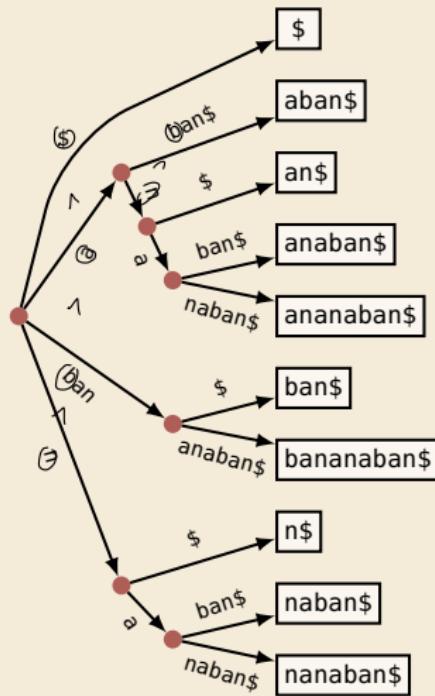
- ◀ construction of suffix trees:
linear time, but significant overhead
- ◀ construction methods fairly complicated
- ◀ many pointers in tree incur large space overhead



13.5 Suffix Arrays

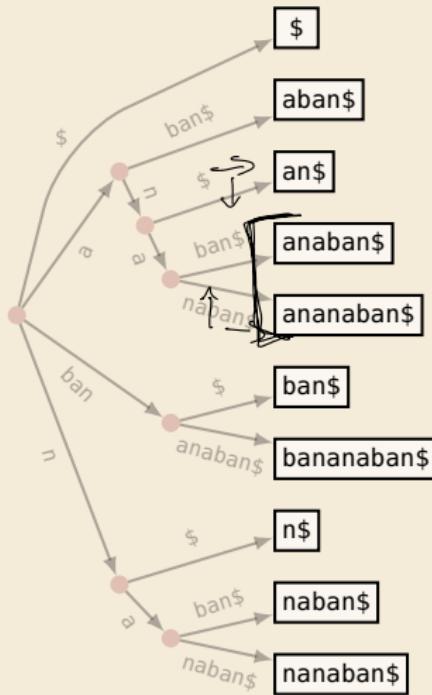
Putting suffix trees on a diet

- ▶ **Observation:** order of leaves in suffix tree
= suffixes lexicographically *sorted*



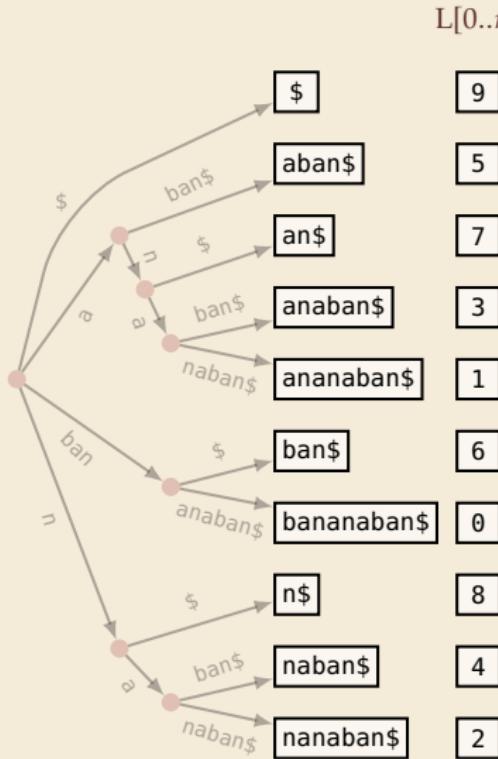
Putting suffix trees on a diet

- ▶ **Observation:** order of leaves in suffix tree
= suffixes lexicographically *sorted*



- ▶ Idea: only store list of leaves $\underline{L[0..n]}$
- ▶ Sufficient to do efficient string matching!
 1. Use binary search for pattern P
 2. check if P is prefix of suffix after position found
- ▶ Example: $P = \underline{\text{ana}}$

Putting suffix trees on a diet



- **Observation:** order of leaves in suffix tree
= suffixes lexicographically *sorted*

- Idea: only store list of leaves $L[0..n]$
- Sufficient to do efficient string matching!
 1. Use binary search for pattern P
 2. check if P is prefix of suffix after position found

- **Example:** $P = \text{ana}$

↔ $L[0..n]$ is called *suffix array*:

$L[r] = (\text{start index of } r\text{th suffix in sorted order})$

- using L , can do string matching with
 $\leq (\lg n + 2) \cdot m$ character comparisons