



# ALGORITHMS OF BIOINFORMATICS

## 6

## Suffix Trees

*11 December 2025*

Prof. Dr. Sebastian Wild



# Outline

## 6 Suffix Trees

- 6.1 Suffix Trees
- 6.2 Applications
- 6.3 Longest Common Extensions
- 6.4 Suffix Arrays
- 6.5 Suffix sorting: Induced sorting and merging
- 6.6 Suffix Sorting: The DC3 Algorithm
- 6.7 The LCP Array
- 6.8 LCP Array Construction



## Context

*We're still working towards practical solutions for the read mapping problem.*

*So far, our preprocessing was mostly getting smart on the **reads/patterns**.*

→ Now preprocess the genome/text.



## 6.1 Suffix Trees

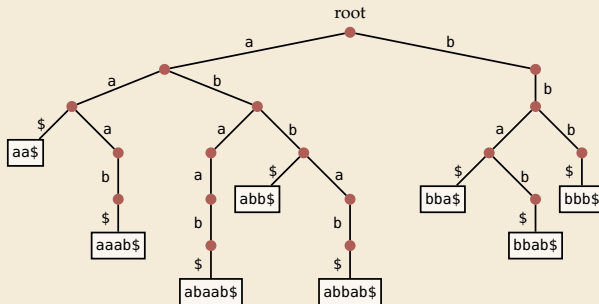


# Recap: Tries

- ▶ efficient dictionary data structure for strings (or for Aho-Corasick automata 😊)
- ▶ name from retrieval, but pronounced “try”
- ▶ tree based on symbol comparisons
- ▶ **Assumption here:** 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\$}

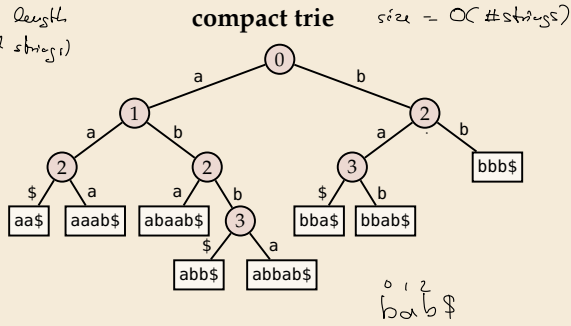
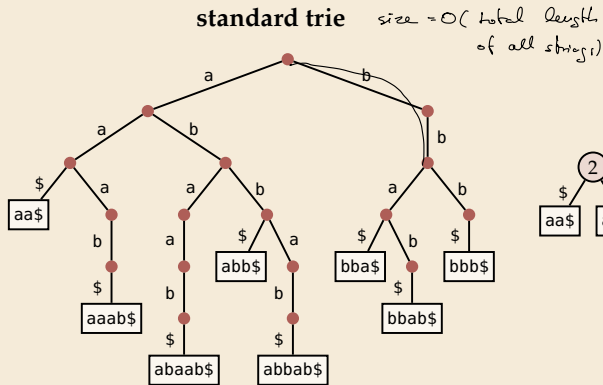




# Compact tries

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

=1 child



- ▶ search gives first character of edge only  $\rightsquigarrow$  must check for match against stored string
- ▶ all nodes  $\geq 2$  children  $\rightsquigarrow \# \text{nodes} \leq \# \text{leaves} = \# \text{strings} \rightsquigarrow$  linear space



# 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



# 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  $T[n] := \$$ )

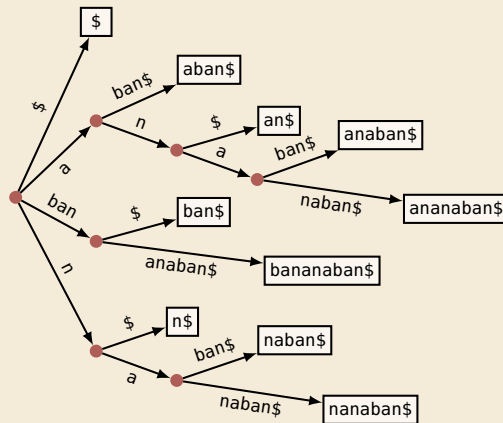
## Example:

$T = \text{bananaban\$}$

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

0	1	2	3	4	5	6	7	8	9
b	a	n	a	n	a	b	a	n	\$

$T =$





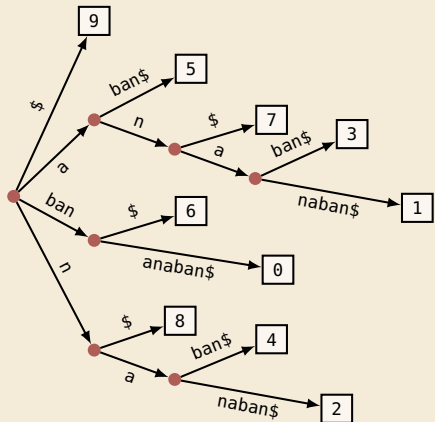
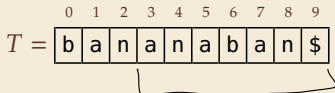
# Suffix trees – Definition

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





# 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)

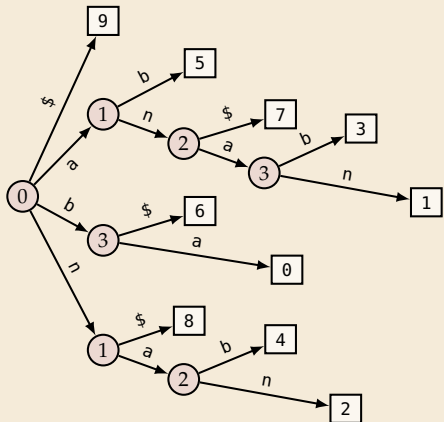
## Example:

$T = \text{bananaban\$}$

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

	0	1	2	3	4	5	6	7	8	9
$T =$	b	a	n	a	n	a	b	a	n	\$

- ▶ 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  $\Theta(n^2)$ .  $\rightsquigarrow$  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  $\Theta(n^2)$ .  $\rightsquigarrow$  not interesting!



same order of growth as reading the text!

**Amazing result:** Can construct the suffix tree of  $T$  in  $\Theta(n)$  time!

- ▶ several fundamentally different methods known
- ▶ started as theoretical breakthrough
- ▶ now routinely used in bioinformatics practice

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



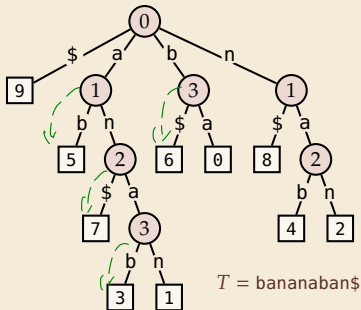
## 6.2 Applications



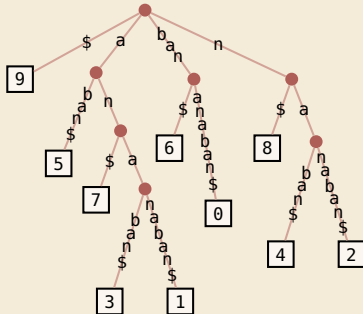
## Applications of suffix trees

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

**Recall:**  $\mathcal{T}$  stored like this:



but think about this:

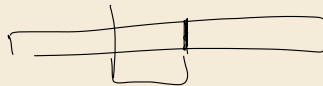


- Moreover: assume internal nodes store pointer to *leftmost leaf in subtree*.
- Notation:  $T_i = T[i..n]$  (including \$)



## 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)

↪  $P$  does not occur in  $T$

2. **we run out of pattern**

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

↪  $P$  occurs at all leaves in subtree of  $v$

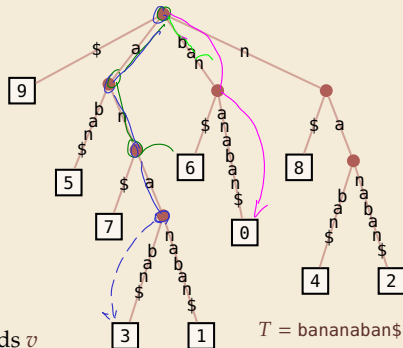
3. **we run out of tree**

reach a leaf  $\ell$  with part of  $P$  left ↪ compare  $P$  to  $\ell$ .



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:

►  $P = \underline{\text{ann}}$

►  $P = \underline{\text{baa}}$

►  $P = \underline{\text{ana}}$

►  $P = \text{ba}$

►  $P = \underline{\text{briar}}$



## 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~~ → 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



How can we efficiently check *all possible substrings*?



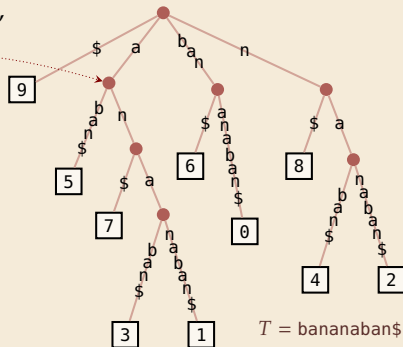
Repeated substrings = shared paths in *suffix tree*



- $T_5 = \text{aban\$}$  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





## Application 2: Longest repeated substring

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



How can we efficiently check *all possible substrings*?

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

