

8

Text Indexing – Searching entire genomes

24 November 2023

Sebastian Wild

Learning Outcomes

- Know and understand methods for text indexing: inverted indices, 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 and LCP arrays.

Unit 8: Text Indexing



Outline

8 Text Indexing

- 8.1 Motivation
- 8.2 Suffix Trees
- 8.3 Applications
- 8.4 Longest Common Extensions
- 8.5 Suffix Arrays
- 8.6 Linear-Time Suffix Sorting: Overview
- 8.7 Linear-Time Suffix Sorting: The DC3 Algorithm
- 8.8 The LCP Array
- 8.9 LCP Array Construction

8.1 Motivation

Text indexing

- ► *Text indexing* (also: *offline text search*):
 - ightharpoonup case of string matching: find P[0..m) in T[0..n)
 - ▶ but with *fixed* text \rightarrow preprocess T (instead of P)
 - \rightarrow expect many queries P, answer them without looking at all of T
 - \leadsto essentially a data structuring problem: "building an *index* of T"

Latin: "one who points out"

- application areas
 - web search engines
 - online dictionaries
 - online encyclopedia
 - ► DNA/RNA data bases
 - ... searching in any collection of text documents (that grows only moderately)

Inverted indices

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

Inverted indices

- ▶ original indices 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:

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

Do you know what a *trie* is?



- 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.



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)
- some character $\notin \Sigma$ ▶ strings of same length 229\$ strings have "end-of-string" marker \$ insert bas root Example: smalle ex. {aa\$,aaab\$,abaab\$,abb\$, abbab\$, bba\$, bbab\$, bbb\$} aa S aa\$ a0a\$ abb\$ bba\$ bbb\$ aaab\$ aas abbab\$ abaab\$

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

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



 $\mathbf{A}) \ \Theta(\log n)$

 \mathbf{F}) $\Theta(\log m)$

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

 $\mathbf{G} \ \Theta(q)$

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

 $oldsymbol{\mathsf{H}} oldsymbol{\Theta}(\log q)$

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

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

 \bullet $\Theta(m)$

 $\Theta(q + \log n)$



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

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



A ⊖(log n)

F∫ Θ(log m

 \mathbf{G} $\Theta(q)$ \checkmark

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

H) Q(log q)

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

 $\Theta(q - \log n)$

E) ⊕(*m*

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



Suppose we have a trie that stores n strings over $\Sigma = \{A, ..., Z\}$. Each stored string consists of m characters. How many **nodes** does the trie have **in total** *in the worst case*?



 $\mathbf{A} \Theta(n)$

 \bigcirc $\Theta(n \log m)$

B $\Theta(n+m)$

 $lackbox{\textbf{E}}$ $\Theta(m)$

 \mathbf{C} $\Theta(n \cdot m)$

 $lackbox{\sf F} \hspace{0.5cm} \Theta(m \log n)$





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

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



A @(11)

Q(n | m)

 \bigcirc $\Theta(n \cdot m) \checkmark$

 $\Theta(n \log m)$

E) ⊕(*m*)

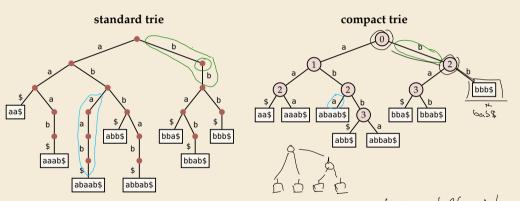


Compact tries

=1 child

bab 3

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



→ searching slightly trickier, but same time complexity as in trie storias cluld pointers ((48)

▶ all nodes \geq 2 children \rightsquigarrow #nodes \leq #leaves = #strings \rightsquigarrow linear space $O(\alpha)$

Tries as inverted index



fast lookup

cannot handle more general queries:

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

Tries as inverted index



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

binary streams

→ need new ideas

8.2 Suffix Trees

Suffix trees – A 'magic' data structure

Appetizer: Longest common substring problem

► Given: strings $S_1, ..., S_k$ Example: S_1 = superiorcalifornializes, S_2 = sealizer

► Goal: find the longest substring that occurs in all *k* strings

Suffix trees – A 'magic' data structure

Appetizer: Longest common substring problem

- ► Given: strings $S_1, ..., S_k$ Example: S_1 = superiorcalifornializes, S_2 = sealizer
- ▶ Goal: find the longest substring that occurs in all k strings \longrightarrow alive



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

Suffix trees – A 'magic' data structure

Appetizer: Longest common substring problem

► Given: strings $S_1, ..., S_k$ Example: S_1 = superiorcalifornializes, S_2 = sealizer

 \blacktriangleright Goal: find the longest substring that occurs in all k strings \leadsto alive



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

Enter: suffix trees

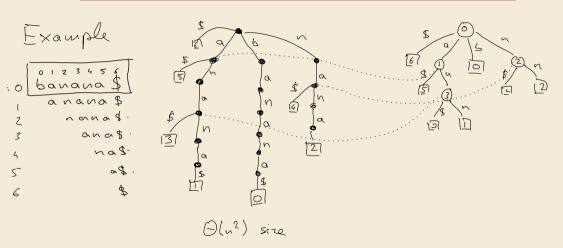
- versatile data structure for index with full-text search
- ▶ linear time (for construction) and linear space
- ▶ allows efficient solutions for many advanced string problems



"Although the longest common substring problem looks trivial now, given our knowledge of suffix trees, it is very interesting to note that in 1970 <u>Don Knuth</u> conjectured that a linear-time algorithm for this problem would be impossible." [Gusfield: Algorithms on Strings, Trees, and Sequences (1997)]

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

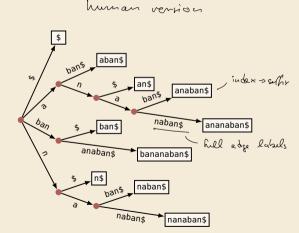
► suffix tree \mathcal{T} for text T = T[0..n) = compact trie of all suffixes of T\$ (set <math>T[n] := \$)



▶ suffix tree \Im for text T = T[0..n) = compact trie of all suffixes of T\$ (set <math>T[n] := \$)

Example:

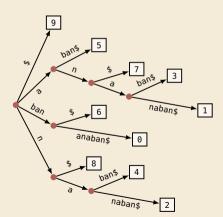
T = bananaban\$



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

Example:

T = bananaban\$

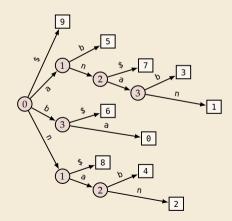


- ► suffix tree T for text T = T[0..n) = compact trie of all suffixes of T\$ (set <math>T[n] := \$)
- except: in leaves, store *start index* (instead of copy of actual string)

Example:

T = bananaban\$

- ▶ 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 $\Theta(n^2)$. \longrightarrow 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 $\Theta(n^2)$. \longrightarrow 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
- 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?

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

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



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

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

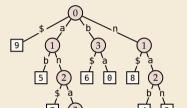


8.3 Applications

Applications of suffix trees

▶ In this section, always assume suffix tree T for T given.

Recall: T stored like this:



but think about this:



▶ Moreover: assume internal nodes store pointer to *leftmost leaf in subtree*.

T = bananaban\$

► Notation: $T_i = T[i..n]$ (including \$)



What does *T*'s suffix tree (on the left) tell you about the question whether T contains the pattern P = ana?

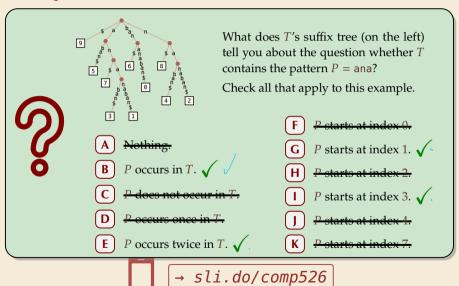
Check all that apply to this example.

- P starts at index 0.
- P starts at index 1.
- P starts at index 2.
- P starts at index 3.
- P starts at index 4.
- P starts at index 7.



- Nothing.
 - P occurs in T.
- P does not occur in T.
- *P* occurs once in *T*.
- *P* occurs twice in *T*.





Application 1: Text Indexing / String Matching

- ▶ we have all suffixes in T!

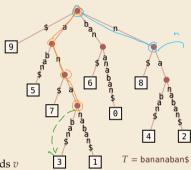
Application 1: Text Indexing / String Matching

- ▶ P occurs in $T \iff P$ is a prefix of a suffix of T
- ▶ we have all suffixes in T!
- \rightsquigarrow (try to) follow path with label P, until
 - 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 ℓ with part of P left \leadsto compare P to ℓ .



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

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



can had all matches by transing subtree

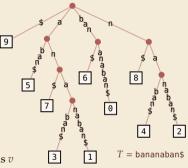
Application 1: Text Indexing / String Matching

- ▶ P occurs in $T \iff P$ is a prefix of a suffix of T
- ▶ we have all suffixes in T!
- \rightsquigarrow (try to) follow path with label P, until
 - 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
 - we run out of tree
 reach a leaf ℓ with part of P left → compare P to ℓ.



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

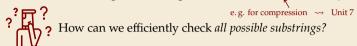
► Finding first match (or NO_MATCH) takes O(|P|) time!



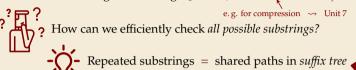
Examples:

- ightharpoonup P = ann
- ightharpoonup P = baa
- ightharpoonup P = ana
- ightharpoonup P = ba
- $P = \underline{b}ri\underline{a}r$

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



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



► T_5 = aban\$ and T_7 = an\$ have longest common prefix 'a' \Rightarrow ∃ internal node with path label 'a' \Rightarrow here single edge, can be longer path



▶ **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*



- $ightharpoonup T_5$ = aban\$ and T_7 = an\$ have longest common prefix 'a'
- → ∃ internal node with path label 'a'

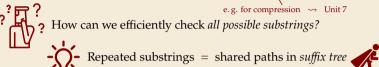
here single edge, can be longer path

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

actually: adjacent leaves



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

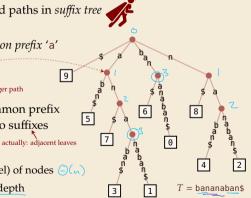


- $ightharpoonup T_5$ = aban\$ and T_7 = an\$ have longest common prefix 'a'
- → ∃ internal node with path label 'a'

here single edge, can be longer path

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

- ► Algorithm:
 - 1. Compute *string depth* (=length of path label) of nodes $\bigcirc(\square)$
 - 2. Find internal nodes with maximal string depth
- ▶ Both can be done in depth-first traversal \rightsquigarrow $\Theta(n)$ time



Generalized suffix trees

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

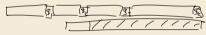
Generalized suffix trees

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

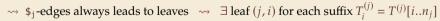
Generalized suffix trees

- ▶ longest *repeated* substring (of one string) feels very similar to longest *common* substring of several strings $T^{(1)}, \ldots, T^{(k)}$ with $T^{(j)} \in \Sigma^{n_j}$
- can we solve that in the same way?
- ightharpoonup 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 $T := T^{(1)} \$_1 T^{(2)} \$_2 \cdots T^{(k)} \$_k$ for k new end-of-word symbols
- ightharpoonup Construct suffix tree T for T





Clicker Question



What is the longest common substring of the strings bcabcac, aabca and bcaa?



→ sli.do/comp526

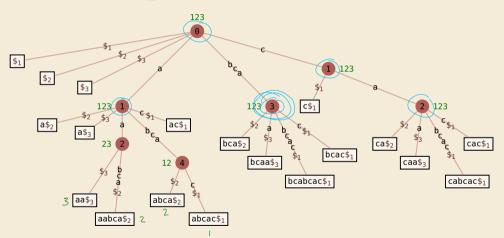
Application 3: Longest common substring

- ▶ With that new idea, we can find longest common substrings:
 - **1.** Compute generalized suffix tree T. \bigcirc (n)
 - **2.** Store with each node the *subset of strings* that contain its path label:
 - 2.1. Traverse 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) $\Diamond(\omega)$
 - **4.** Report deepest node (by string depth) whose subset is $\{1, \ldots, k\}$.
- ▶ Each step takes time $\Theta(n)$ for $n = n_1 + \cdots + n_k$ the total length of all texts.

[&]quot;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

$$T^{(1)} = \text{bcabcac}, \quad T^{(2)} = \text{aab}\underline{\text{ca}}, \quad T^{(3)} = \text{bcaa}$$



8.4 Longest Common Extensions

Application 4: Longest Common Extensions

▶ We implicitly used a special case of a more general, versatile idea:

The *longest common extension (LCE)* data structure:

► Given: String
$$T[0..n)$$

Coal: Answer LCE queries, i. e.,
given positions i, j in T ,
how far can we read the same text from there?
formally: LCE $(i, j) = \max\{\ell : T[i..i + \ell) = T[j..j + \ell)\}$

Application 4: Longest Common Extensions

▶ We implicitly used a special case of a more general, versatile idea:

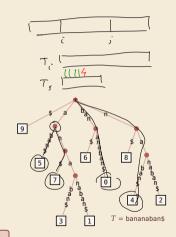
The *longest common extension (LCE)* data structure:

- ▶ **Given:** String T[0..n)
- ► **Goal:** Answer LCE queries, i. e., given positions *i*, *j* in *T*, how far can we read the same text from there?
 - formally: LCE $(i, j) = \max\{\ell : T[i..i + \ell) = T[j..j + \ell)\}$
- \rightsquigarrow use suffix tree of T!

(length of) longest common prefix of ith and jth suffix

In \mathfrak{T} : LCE(i,j) = LCP (T_i,T_j) \longrightarrow same thing, different name! = string depth of lowest common ancester (LCA) of leaves i and j

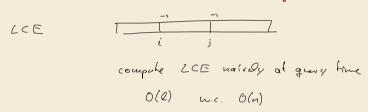
▶ in short: $LCE(i, j) = LCP(T_i, T_j) = stringDepth(LCA(i, j))$



Efficient LCA

How to find lowest common ancestors?

- ► Could walk up the tree to find LCA \rightsquigarrow $\Theta(n)$ worst case
- ► Could store all LCAs in big table \rightsquigarrow $\Theta(n^2)$ space and preprocessing \bigcirc



Efficient LCA

How to find lowest common ancestors?

- ► Could walk up the tree to find LCA \rightsquigarrow $\Theta(n)$ worst case
- ► Could store all LCAs in big table \longrightarrow $\Theta(n^2)$ space and preprocessing



Amazing result: Can compute data structure in $\Theta(n)$ time and space that finds any LCA is **constant(!) time**.

- ► a bit tricky to understand
- but a theoretical breakthrough
- and useful in practice

and suffix tree construction inside . . .

 \rightarrow for now, use O(1) LCA as black box.



~ Duit 9

 \rightarrow After linear preprocessing (time & space), we can find LCEs in O(1) time.

Application 5: Approximate matching

k-mismatch matching:

- ▶ **Input:** text T[0..n), pattern P[0..m), $k \in [0..m)$
- ► Output: "Hamming distance ≤ k"
 - \blacktriangleright smallest *i* so that T[i..i + m) are *P* differ in at most *k* characters
 - ightharpoonup or NO_MATCH if there is no such i
- → searching with typos
- ► Adapted brute-force algorithm \rightsquigarrow $O(n \cdot m)$

Application 5: Approximate matching

k-mismatch matching:

- ▶ **Input:** text T[0..n), pattern P[0..m), $k \in [0..m)$
- ► Output:

```
"Hamming distance \leq k"
```

- \blacktriangleright smallest *i* so that T[i..i + m) are *P* differ in at most *k* characters
- ightharpoonup or NO MATCH if there is no such i
- → searching with typos



- ▶ Adapted brute-force algorithm \rightsquigarrow $O(n \cdot m)$
- Assume longest common extensions in $T $_1P $_2$ can be found in O(1)
 - → generalized suffix tree T has been built
 - → string depths of all internal nodes have been computed
 - → constant-time LCA data structure for T has been built

Kangaroo Algorithm for approximate matching



```
procedure kMismatch(T[0..n-1], P[0..m-1])

// build LCE data structure

for i := 0, ..., n-m-1 do

mismatches := 0; t := i; p := 0

while mismatches \le k \land p < m do

\ell := \underbrace{LCE}(t, p) / jump \ over \ matching \ part

t := \widehat{t + \ell} + 1; \ p := p + \ell + 1

mismatches := \min mismatches + 1

if p == m then

return i
```

- ▶ Analysis: $\Theta(n+m)$ preprocessing + $O(n \cdot k)$ matching
- \rightsquigarrow very efficient for small k
- ▶ State of the art
 - $ightharpoonup O(n^{\frac{k^2 \log k}{m}})$ possible with complicated algorithms
- # exam

ightharpoonup extensions for edit distance $\leq k$ possible

Application 6: Matching with wildcards

► Allow a wildcard character in pattern stands for arbitrary (single) character unit* P in_unit5_we_will T

▶ similar algorithm as for *k*-mismatch \rightsquigarrow $O(n \cdot k + m)$ when *P* has *k* wildcards

Application 6: Matching with wildcards

- ► Allow a wildcard character in pattern stands for arbitrary (single) character
- $\begin{array}{ccc} & & & & P \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & \\ & & & \\ &$
- ▶ similar algorithm as for *k*-mismatch \rightsquigarrow $O(n \cdot k + m)$ when *P* has *k* wildcards

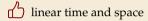
* * *

Many more applications, in particular for problems on biological sequences

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

Suffix trees – Discussion

► Suffix trees were a threshold invention



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

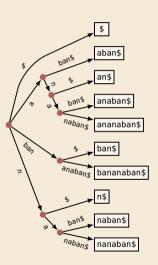
many pointers in tree incur large space overhead



node es arrey dild[0.5) es 5 space overhed
. BST ones chied labels - ls 5 time overhed

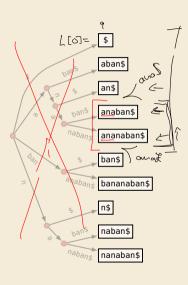
8.5 Suffix Arrays

Putting suffix trees on a diet



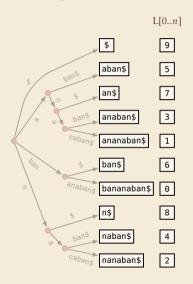
► **Observation:** order of leaves in suffix tree = suffixes lexicographically *sorted*

Putting suffix trees on a diet



- ► **Observation:** order of leaves in suffix tree = suffixes lexicographically *sorted*
- ▶ Idea: only store list of leaves L[0..n]
- Enough to do efficient string matching!
 - 1. Use binary search for pattern P and & and &
 - 2. check if *P* is prefix of suffix after position found
- ightharpoonup **Example:** P = ana

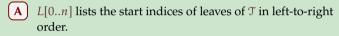
Putting suffix trees on a diet



- ► **Observation:** order of leaves in suffix tree = suffixes lexicographically *sorted*
- ▶ Idea: only store list of leaves L[0..n]
- Enough 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 = ana
- \rightsquigarrow L[0..n] is called *suffix array*:
 - L[r] = (start index of) rth suffix in sorted order
- ▶ using L, can do string matching with $\leq (\lg n + 2) \cdot m$ character comparisons

Clicker Question

Check all correct statements about *suffix array* L[0..n] and *suffix tree* \mathfrak{T} of text T[0..n) (for $\sigma = O(1)$)



- **B** T[L[r]..n] is the path label in \mathcal{T} to the leaf storing r.
- C T[L[r]..n] is the path label to the rth leaf in T.
- **D** $T_{L[r]}$ is the rth smallest suffix of T (lexicographic order).
- **E** In terms of Θ-classes, \mathcal{T} needs more space than L.
- $oldsymbol{\mathsf{F}}$ L (and T) suffice to solve the text indexing problem.



→ sli.do/comp526



Clicker Question

Check all correct statements about *suffix array* L[0..n] and *suffix tree* \Im of text T[0..n) (for $\sigma = O(1)$)

- A L[0..n] lists the start indices of leaves of T in left-to-right order.
- B T[L[r]..n] is the path label in T to the leaf storing r.
- C T[L[r]..n] is the path label to the rth leaf in \mathcal{T} . \checkmark
- D $T_{L[r]}$ is the rth smallest suffix of T (lexicographic order). \checkmark
- E In terms of Θ-classes, T needs more space than L.
- $oxed{\mathsf{F}}$ L (and T) suffice to solve the text indexing problem. \checkmark



 \rightarrow sli.do/comp526



Suffix arrays – Construction

How to compute L[0..n]?

- ▶ from suffix tree
 - ▶ possible with traversal . . .
 - but we are trying to avoid constructing suffix trees!
- ▶ sorting the suffixes of *T* using general purpose sorting method
 - trivial to code!

1 = aaaaaaa

- **b** but: comparing two suffixes can take $\Theta(n)$ character comparisons
- \bigcap $\Theta(n^2 \log n)$ time in worst case

Suffix arrays – Construction

How to compute L[0..n]?

- ▶ from suffix tree
 - ▶ possible with traversal . . .
 - $\hfill \Box$ but we are trying to avoid constructing suffix trees!
- ▶ sorting the suffixes of *T* using general purpose sorting method
 - trivial to code!
 - **b** but: comparing two suffixes can take $\Theta(n)$ character comparisons
 - \bigcap $\Theta(n^2 \log n)$ time in worst case
- ▶ We can do better!

Clicker Question

What is the relation between suffix array L[0..n] and BWT B[0..n] of a string T[0..n)\$?



- $oldsymbol{\mathsf{A}}$ L can be very easily computed from B and T
- f B B can be very easily computed from L and T
- C Both A and B
- **D** Neither A nor B



→ sli.do/comp526

Clicker Question

What is the relation between suffix array L[0..n] and BWT B[0..n] of a string T[0..n)\$?



- (A) Lean be very easily computed from Beand T
- **B** B can be very easily computed from L and T
- Both A and B
- Neither A nor B



→ sli.do/comp526

Digression: Recall BWT

Burrows-Wheeler Transform

- **1.** Take all cyclic shifts of *S*
- 2. Sort cyclic shifts
- 3. Extract last column

S = alf_eats_alfalfa\$
B = asff\$f.e.lllaaata

alf_eats_alfalfa\$ lf..eats..alfalfa\$a f_eats_alfalfa\$al _eats_alfalfa\$alf eats_alfalfa\$alf ats alfalfa\$alf e ts_alfalfa\$alf.ea s..alfalfa\$alf..eat _alfalfa\$alf_eats alfalfa\$alf,eats, lfalfa\$alf..eats..a falfa\$alf.eats.al alfa\$alf_eats_alf lfa\$alf_eats_alfa fa\$alf_eats_alfal a\$alf, eats, alfalf \$alf, eats, alfalfa

\$alf,eats,alfalfa ..alfalfa\$alf..eats _eats_alfalfa\$alf a\$alf_eats_alfalf alf_eats_alfalfa\$ alfa\$alf_eats_alf alfalfa\$alf.eats.. ats alfalfa\$alf_e eats_alfalfa\$alf. f.,eats,,alfalfa\$a<mark>t</mark> fa\$alf,eats,alfal falfa\$alf_eats_al lf.eats.alfalfa\$a lfa\$alf_eats_alfa lfalfa\$alf_eats_a s..alfalfa\$alf_eat ts..alfalfa\$alf..ea

 \longrightarrow

sort

Digression: Computing the BWT

How can we compute the BWT of a text efficiently?

Digression: Computing the BWT

How can we compute the BWT of a text efficiently?

- ightharpoonup cyclic shifts S = suffixes of S
 - comparing cyclic shifts stops at first \$
 - ► for comparisons, anything after \$ irrelevant!
- ► BWT is essentially suffix sorting!
 - ► B[i] = S[L[i] 1]
 - ▶ where L[i] = 0, B[i] = \$
- \rightsquigarrow Can compute B in O(n) time from L

```
alf, eats, alfalfa$
lf.eats.alfalfa$a
f, eats, alfalfa$al
_eats_alfalfa$alf
eats, alfalfa$alf...
ats.alfalfa$alf.e
ts.,alfalfa$alf.,ea
s..alfalfa$alf..eat
..alfalfa$alf..eats
alfalfa$alf..eats..
lfalfa$alf,.eats,.a
falfa$alf_eats_al
alfa$alf,.eats,.alf
lfa$alf..eats..alfa
fa$alf, eats, alfal
a$alf..eats..alfalf
$alf, eats, alfalfa
```

```
\downarrow L[r]
   $alf,_eats_alfalfa
   ..alfalfa$alf..eats
   .,eats,,alfalfa$alf
   a$alf_eats_alfalf
   alf.eats.alfalfa$
   alfa$alf,eats,alf
   alfalfa$alf..eats...
   ats.alfalfa$alf.e
                         5
   eats, alfalfa$alf...
   f..eats..alfalfa$al
10 fa$alf,eats,alfal
   falfa$alf,,eats,,al
12 lf.eats.alfalfa$a
  lfa$alf..eats..alfa
                        13
   lfalfa$alf,.eats..a
                        10
   s..alfalfa$alf,eat
   ts.alfalfa$alf.ea
```

Fat-pivot radix quicksort – Example

she **s**ells **s**eashells by the sea **s**hore the **s**hells she **s**ells are **s**urely **s**eashells

she

sells

seashells

by

the

sea

shore

the

shells

she

sells

are

surely

seashells

s he	b y
s ells	a re
s eashells	s he
b y	s e lls
the	s e ashells
s ea	s e a
s hore	shore
the	s hells
s hells	s he
s he	s e lls
s ells	surely
a re	s e ashells
s urely	the
s eashells	t he

she	b y	are
s ells	are	by
s eashells	s h e	
b y	s e lls	
the	s e ashells	
sea	s e a	
shore	s h ore	
the	shells	
s hells	s he	
she	s e lls	
s ells	surely	
a re	seashells	
s urely	the	
seashells	the	

she	b y	are
s ells	a re	by
s eashells	she	s e lls
b y	s e lls	s e ashells
the	s e ashells	s e a
sea	s e a	s e lls
shore	s h ore	s e ashells
the	shells	sh e
s hells	s he	shore
she	s e lls	sh e lls
s ells	surely	sh e
are	s e ashells	surely
surely	the	
seashells	the	

she	by	are
s ells	a re	by
s eashells	she	sells
b y	s e lls	s e ashells
the	s e ashells	s e a
sea	s e a	s e lls
shore	shore	s e ashells
the	shells	she
s hells	she	shore
she	s e lls	sh ells
s ells	surely	she
are	s e ashells	surely
s urely	the	the
seashells	the	the

s he	b y	are			
s ells	a re	by			
s eashells	she	sells		se lls	
b y	s e lls	s e ashells		se a shells	
t he	s e ashells	s ea		sea	
s ea	s e a	s ells		se lls	
s hore	s h ore	s e ashells		se a shells	
t he	s hells	she	Ì		
s hells	s he	sh ore			
s he	s e lls	sh ells			
s ells	surely	sh e			
a re	s e ashells	surely			
s urely	the	the			
s eashells	t he	the			

s he	b y	are		
s ells	a re	by		
s eashells	she	sells	se l ls	Ì
b y	s e lls	s e ashells	se a shells	
t he	s e ashells	s e a	se a	
s ea	s ea	s e lls	se lls	
s hore	s h ore	s e ashells	se a shells	
t he	s h ells	she	she \$	1
s hells	s he	shore	she lls	
s he	s e lls	sh e lls	she \$	
s ells	surely	sh e	shore	
a re	s e ashells	surely		
s urely	the	t h e		
s eashells	the	t h e		

s he	b y	are		
s ells	a re	by		
s eashells	she	sells	se l ls	
b y	s e lls	s e ashells	se a shells	
t he	s e ashells	s ea	sea	
s ea	s ea	s e lls	se lls	
s hore	s h ore	s e ashells	se a shells	
the	s hells	she	she \$	
s hells	s he	shore	she lls	
s he	s e lls	sh ells	she \$	
s ells	surely	sh e	shore	
a re	s e ashells	surely		
s urely	the	the	the	
s eashells	the	t h e	the	

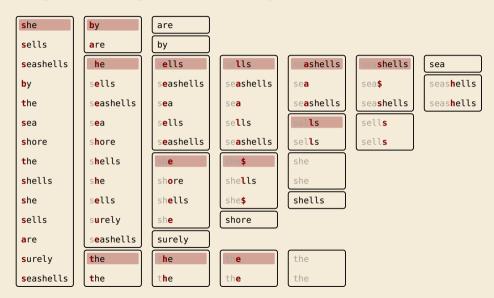
she	by	are		
sells	a re	by		
s eashells	she	ells	sells	seashells
b y	s e lls	seashells	se a shells	sea
the	s e ashells	s e a	sea	se a shells
sea	s e a	s e lls	se lls	sel ls
shore	s h ore	seashells	se a shells	sel ls
the	s h ells	she	she \$	
s hells	s he	shore	she lls	
s he	s e lls	sh e lls	she \$	
s ells	surely	sh e	shore	
are	s e ashells	surely		
s urely	the	the	the	
s eashells	t he	t h e	the	

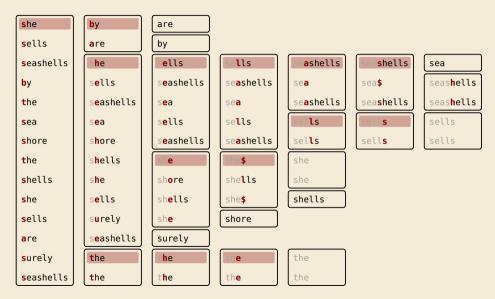
she	b y	are		
sells	a re	by		
s eashells	she	sells	sells	seashells
b y	s e lls	s e ashells	se a shells	sea
the	s e ashells	s e a	sea	se a shells
sea	s e a	s e lls	se lls	sel ls
shore	s h ore	s e ashells	se a shells	sel ls
the	shells	she	she\$	she
s hells	s he	shore	she lls	she
she	s e lls	sh e lls	she \$	shells
sells	surely	sh e	shore	
are	s e ashells	surely		
surely	the	the	the	
s eashells	the	t h e	the	

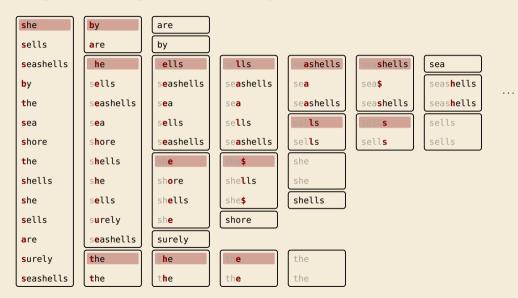
she	b y	are		
s ells	a re	by		
s eashells	she	sells	sells	se a shells
b y	s ells	s e ashells	se a shells	sea
the	s e ashells	s ea	se a	se a shells
s ea	s ea	s e lls	se l ls	sel ls
s hore	s hore	s e ashells	se a shells	sel ls
the	s hells	she	she\$	she
s hells	s he	shore	she lls	she
s he	s ells	sh ells	she \$	shells
s ells	surely	sh e	shore	
a re	s e ashells	surely		
s urely	the	the	the	the
s eashells	the	t h e	the	the

she	b y	are			
s ells	a re	by			
s eashells	she	sells	sells	seashells	sea s hells
b y	s e lls	s e ashells	se a shells	se a	sea \$
t he	s e ashells	s ea	se a	se a shells	sea shells
sea	s ea	s e lls	se lls	sel ls	
s hore	s h ore	s e ashells	se a shells	sel ls	
t he	s hells	she	she \$	she	
s hells	s he	shore	she lls	she	
s he	s e lls	sh ells	she \$	shells	
s ells	surely	sh e	shore		
a re	s e ashells	surely			
s urely	the	the	the	the	
s eashells	t he	t h e	the	the	

she	by	are			
s ells	a re	by			
s eashells	she	sells	sells	sashells	sea s hells
b y	s e lls	s e ashells	se a shells	se a	sea \$
the	s e ashells	s e a	se a	se a shells	sea s hells
sea	s ea	s e lls	se l ls	sells	sells
shore	s h ore	s e ashells	se a shells	sel ls	sell s
the	shells	she	she\$	she	
s hells	s he	shore	she lls	she	
s he	s e lls	sh e lls	she \$	shells	
s ells	surely	sh e	shore		
are	s e ashells	surely			
surely	the	the	the	the	
s eashells	the	t h e	the	the	







Fat-pivot radix quicksort

details in §5.1 of Sedgewick, Wayne Algorithms 4th ed. (2011), Pearson

- **partition** based on *d*th character only (initially d = 0)
- → 3 segments: smaller, equal, or larger than *d*th symbol of pivot
- recurse on smaller and large with same d, on equal with d + 1
 - $\rightsquigarrow \ never \ compare \ equal \ prefixes \ twice$

Fat-pivot radix quicksort

details in §5.1 of Sedgewick, Wayne Algorithms 4th ed. (2011), Pearson

- **partition** based on *d*th character only (initially d = 0)
- \rightarrow 3 segments: smaller, equal, or larger than dth symbol of pivot
- recurse on smaller and large with same d, on equal with d+1
 - → never compare equal prefixes twice

for random strings

- \rightarrow can show: $\sim 2 \ln(2) \cdot n \lg n \approx 1.39 n \lg n$ character comparisons on average
- simple to code
- efficient for sorting many lists of strings

random string

• fat-pivot radix quicksort finds suffix array in $O(n \log n)$ expected time

Fat-pivot radix quicksort

details in §5.1 of Sedgewick, Wayne Algorithms 4th ed. (2011), Pearson

- **partition** based on *d*th character only (initially d = 0)
- \rightarrow 3 segments: smaller, equal, or larger than dth symbol of pivot
- recurse on smaller and large with same d, on equal with d + 1
 - → never compare equal prefixes twice

for random strings

- \rightarrow can show: $\sim 2 \ln(2) \cdot n \lg n \approx 1.39 n \lg n$ character comparisons on average
- simple to code
- efficient for sorting many lists of strings

random string

▶ fat-pivot radix quicksort finds suffix array in $O(n \log n)$ expected time

but we can do O(n) time worst case!