6

Text Indexing – Searching whole genomes

7 March 2022

Sebastian Wild

Learning Outcomes

- 1. 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 6: Text Indexing



Outline

6 Text Indexing

- 6.1 Motivation
- 6.2 Suffix Trees
- 6.3 Applications
- 6.4 Longest Common Extensions
- 6.5 Suffix Arrays
- 6.6 Linear-Time Suffix Sorting: Overview
- 6.7 Linear-Time Suffix Sorting: The DC3 Algorithm
- 6.8 The LCP Array
- 6.9 LCP Array Construction

6.1 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 \leadsto preprocess T (instead of P)
 - \rightarrow expect many queries P, answer them without looking at all of T
 - → 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

- same as "indexes"
- ▶ 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

- \triangleright original indices in books: list of (key) words \mapsto 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 whitespace
 - actual implementations typically support stemming of words

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.

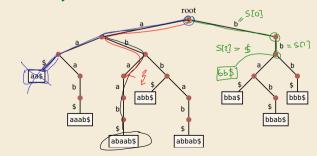
Tries

- 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
- strings have "end-of-string" marker \$
- Example: $Z = \{a,b\}$ $\{\underline{aa},\underline{aa}$ ab\$, abaab\$, abb\$, abbab\$, bba\$, bbab\$, bbb\$}

insurt S=b5\$ aba\$\$ ST



does the string occur?

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



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

 $\Theta(\log(nm))$

 \bigcirc $\Theta(q)$

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

 $(\mathbf{H}) \ \Theta(\log q)$

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

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

 \mathbf{E} $\Theta(m)$

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

sli.do/comp526

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**?



- $\mathbf{B} \quad \Theta(\log(nm)) \qquad \qquad \mathbf{G} \quad \Theta(q) \quad \checkmark$
- $\mathsf{C} \hspace{.1cm} \hspace{.1cm} \hspace{.1cm} \Theta(m \hspace{.1cm} \hspace{.1cm} \log n) \hspace{1cm} \hspace{.1cm} (\mathsf{H}) \hspace{.1cm} \hspace{.1cm} \Theta(\log q)$

sli.do/comp526

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*?



 $oldsymbol{\mathsf{A}} \hspace{0.1cm} \Theta(n)$

 $\Theta(n+m)$

 $oldsymbol{\mathsf{E}} oldsymbol{\Theta}(m)$

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

 $\Theta(m \log n)$



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)

 $\Theta(n \perp m)$

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

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

E) ⊕(m)

 $\Theta(m \log n)$

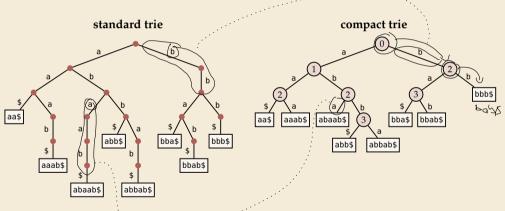
Compact tries

=1 child



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
- ▶ all nodes ≥ 2 children → #nodes ≤ #leaves = #strings → linear space



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

binary streams

→ need new ideas

6.2 Suffix Trees

Suffix trees – A 'magic' data structure

Appetizer: Longest common substring problem

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



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

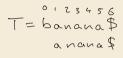
Enter: suffix trees

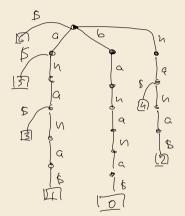
- versatile data structure for index with full-text search
- ▶ linear time (for construction) and linear space
- allows efficient solutions for many advanced string problems

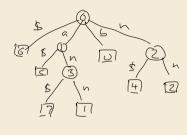


"Although the longest common substring problem looks trivial now, given our knowledge of suffix trees, it is very interesting to note that in 1970 Don Knuth conjectured that a linear-time algorithm for this problem would be impossible." [Gusfield: Algorithms on Strings, Trees, and Sequences (1997)]

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



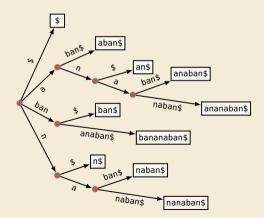




▶ suffix tree \mathcal{T} for text T = T[0..n) = compact trie of all suffixes of T\$ (set <math>T[n] := \$)

Example:

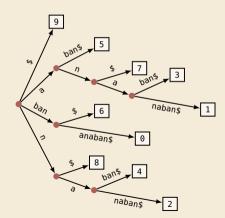
T = bananaban\$



- suffix tree \Im 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\$

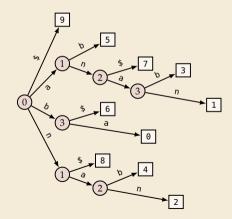


- ▶ suffix tree \Im for text T = T[0..n) = compact trie of all suffixes of T\$ (set <math>T[n] := \$)
- except: in leaves, store *start index* (instead of copy of actual string)

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).

- $oxed{A}$ 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).
- H) Thas n leaves.

%

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

- A We require T to end with \$. \checkmark
- B The size of T can be $\Omega(n^2)$ in the worst case.
- T is a standard trie of all suffixes of T\$.
- **D** T is a compact trie of all suffixes of T\$. \checkmark
- 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). \checkmark
- G T can be computed in O(n) time (worst case). \checkmark
- H Thas n leaves.

6.3 Applications

Applications of suffix trees

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

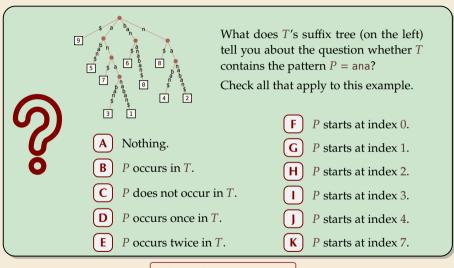
Recall: T stored like this:

9 1 3 1 1 5 2 6 0 8 2 5 7 3 4 2 7 3 T = bananaban\$

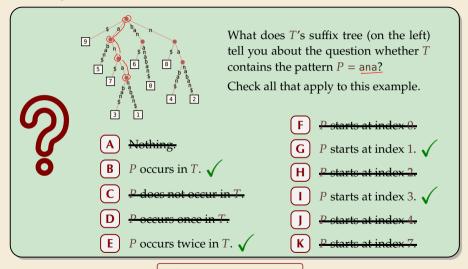
but think about this:



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



sli.do/comp526



sli.do/comp526

Application 1: Text Indexing / String Matching

- ▶ P occurs in T \iff P is a prefix of a suffix of T
- \blacktriangleright we have all suffixes in T!

Application 1: Text Indexing / String Matching

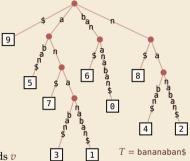
- ▶ P occurs in $T \iff P$ is a prefix of a suffix of T
- ▶ we have all suffixes in T!
- \rightsquigarrow (try to) follow path with label P, until
 - we get stuck
 at internal node (no node with next character of P)
 or inside edge (mismatch of next characters)
 → P does not occur in T
 - 2. we run out of pattern

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

 \rightarrow *P* occurs at all leaves in subtree of *v*

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!

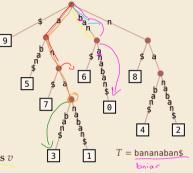


follow leftwost-lost pointe

Application 1: Text Indexing / String Matching

- ightharpoonup 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
 - 1. we get stuck

 at internal node (no node with next character of P)
 or inside edge (mismatch of next characters)
 - \rightarrow P does not occur in T
 - 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 $\$ \notin \Sigma$, 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 \$
- P = ana(
 - \triangleright P = briar



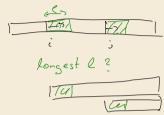
Application 2: Longest repeated substring

▶ **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 \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



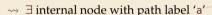
How can we efficiently check all possible substrings?



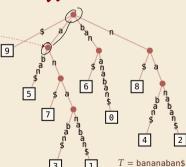
Repeated substrings = shared paths in *suffix tree*



 $ightharpoonup T_5$ = aban\$ and T_7 = an\$ have longest common prefix 'a'

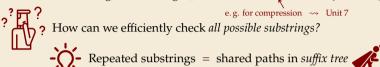


ere single edge, can be longer path



Application 2: Longest repeated substring

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



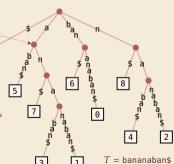


- $ightharpoonup T_5 = aban\$$ and $T_7 = an\$$ have longest common prefix 'a'
- → ∃ internal node with path label 'a'

here single edge, can be longer path

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

actually: adjacent leaves



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*

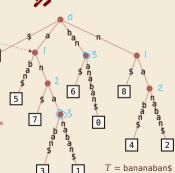


- $ightharpoonup T_5 = aban\$$ and $T_7 = an\$$ have longest common prefix 'a'
- → ∃ internal node with path label 'a'

here single edge, can be longer path

- longest repeated substring = longest common prefix (LCP) of two suffixes
 - actually: adjacent leaves

- ► Algorithm:
 - 1. Compute string depth (=length of path label) of nodes
 - 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 $T^{(1)}, \ldots, T^{(k)}$ 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

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
- \Rightarrow \$j-edges always leads to leaves \Rightarrow \exists leaf (j,i) for each suffix $T_i^{(j)} = T^{(j)}[i..n_j]$



Clicker Question



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

sli.do/comp526

Application 3: Longest common substring

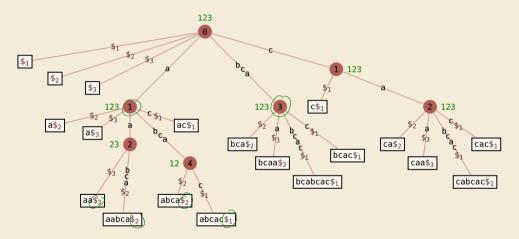
- ▶ With that new idea, we can find longest common substrings:
 - **1.** Compute generalized suffix tree T.
 - **2.** Store with each node the *subset of strings* that contain its path label:
 - 2.1. Traverse 𝒯 bottom-up.
 - **2.2.** For a leaf (j, i), the subset is $\{j\}$.

- O(u) time
- 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 $\{1, \ldots, k\}$.
- ▶ 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

 $T^{(1)}=$ bcabcac, $T^{(2)}=$ aabca, $T^{(3)}=$ bcaa



6.4 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:

MA WO

- ▶ **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{ $\ell : T[i..i + \ell) = T[j..j + \ell)$ }

Application 4: Longest Common Extensions

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

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

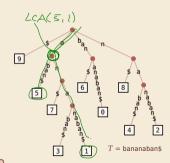
- ▶ **Given:** String T[0..n)
- ► **Goal:** Answer LCE queries, i. e., given positions *i*, *j* in *T*, how far can we read the same text from there?
 - formally: LCE $(i, j) = \max\{\ell : T[i..i + \ell) = T[j..j + \ell)\}$
- \rightsquigarrow use suffix tree of T!

(length of) longest common prefix

of *i*th and *j*th suffix

- ► In \mathfrak{T} : LCE $(i, j) = \text{LCP}(T_i, T_j) \rightsquigarrow \text{ same thing, different name!}$ = string depth of lowest common ancester (LCA) of leaves i and j
- ▶ in short: $LCE(i, j) = LCP(T_i, T_j) = stringDepth(LCA(i, j))$

LCP (Ts, T,)



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 \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



- \rightarrow for now, use O(1) LCA as black box.
- \rightarrow After linear preprocessing (time & space), we can find <u>LCEs</u> 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 $O(n \cdot m)$



compare w/ exact striss
matching O(n+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 $T \$_1 P \$_2$ 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

Clicker Question



What is the Hamming distance between heart and beard?

2

sli.do/comp526

Kangaroo Algorithm for approximate matching

```
STE
```

```
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 \leq k \wedge p < m do

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

t := t + \ell + 1; p := p + \ell + 1

mismatches := mismatches + 1

if p == m then

return i
```

- ▶ **Analysis:** $\Theta(n+m)$ preprocessing + $O(n \cdot k)$ matching
- ► State of the art
 - $O(n^{\frac{k^2 \log k}{m}})$ possible with complicated algorithms
 - \blacktriangleright 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

 unit*

 in_unit5_uwe_uwill

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

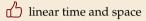
* * *

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

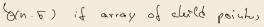


Suffix trees – Discussion

- 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



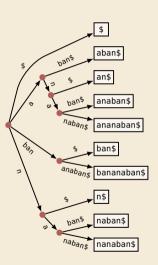
O(n 9,8) if we use dichocain, but the search



contine 11:52

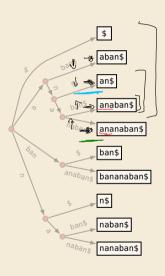
6.5 Suffix Arrays

Putting suffix trees on a diet



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

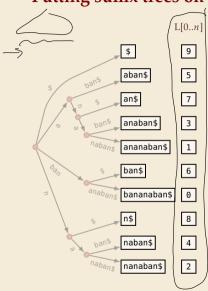
Putting suffix trees on a diet



- ► **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*
 - 2. check if *P* is prefix of suffix after position found



Putting suffix trees on a diet



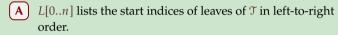
- ► **Observation:** order of leaves in suffix tree = suffixes lexicographically *sorted*
- ▶ Idea: only store list of leaves L[0..n]
- Enough to do efficient string matching!
 - **1.** Use binary search for pattern *P*
 - **2.** check if *P* is prefix of suffix after position found
- **Example:** P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

$$L[r] =$$
(start index of) r th 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 $\sigma = O(1)$)



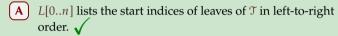


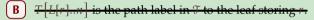
- **B** T[L[r]..n] is the path label in \mathcal{T} to the leaf storing r.
- T[L[r]..n] is the path label to the rth leaf in T.
- **D** $T_{L[r]}$ is the rth smallest suffix of T (lexicographic order).
- **E** In terms of Θ -classes, \mathcal{T} needs more space than L.
- $oldsymbol{\mathsf{F}}$ L (and 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* \mathfrak{T} of text T[0..n) (for $\underline{\sigma} = O(1)$)





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 ⊕ classes, T needs more space than L.

F L (and T) suffice to solve the text indexing problem. \checkmark

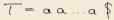


sli.do/comp526

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!



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

- 1) better string sorting

she **s**ells seashells by the sea shore the **s**hells she **s**ells **a**re surely **s**eashells

she

sells

seashells

by

the

sea

shore

the

shells

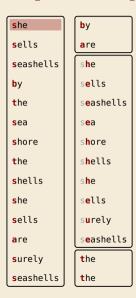
she

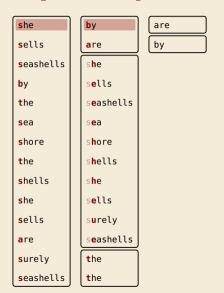
sells

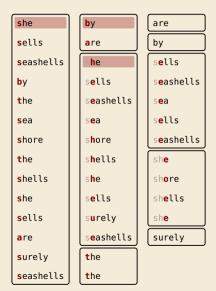
are

surely

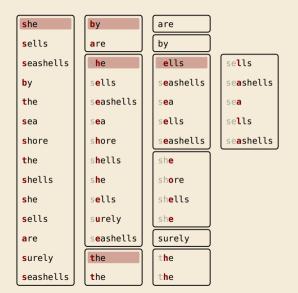
seashells

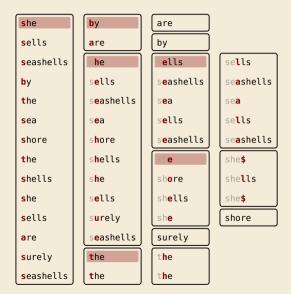




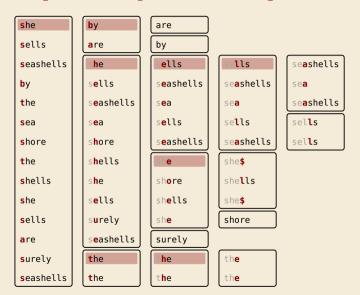


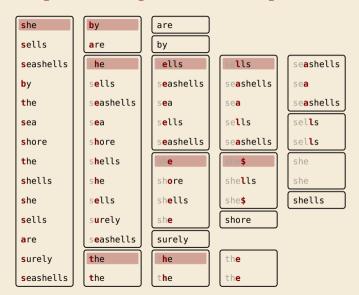
she	b y	are
s ells	a re	by
s eashells	she	s e lls
b y	s e lls	s e ashells
the	s e ashells	s ea
sea	s e a	s e lls
shore	s h ore	s e ashells
the	shells	she
s hells	s he	shore
s he	s e lls	sh e lls
s ells	surely	sh e
are	s e ashells	surely
surely	the	the
seashells	the	the



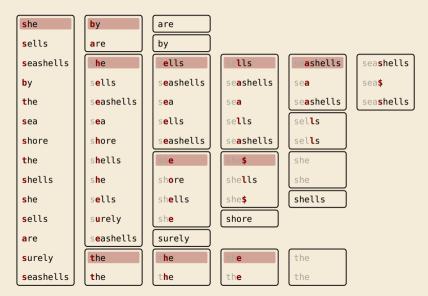


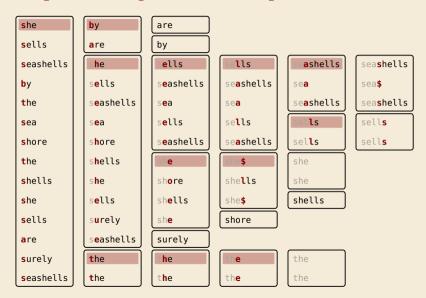
she	by	are	
sells	a re	by	
seashells	she	sells	se l ls
b y	s e lls	s e ashells	se a shells
the	s e ashells	s e a	se a
sea	s e a	s e lls	se l ls
shore	s h ore	s e ashells	se a shells
the	shells	she	she \$
s hells	s he	shore	she lls
she	s e lls	sh e lls	she \$
s ells	s u rely	sh e	shore
are	s e ashells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

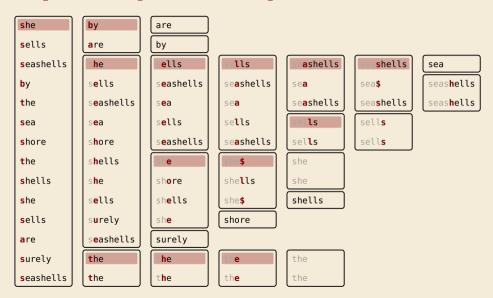


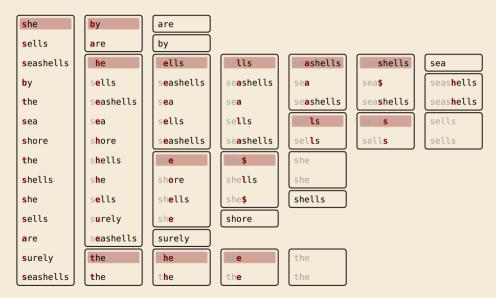


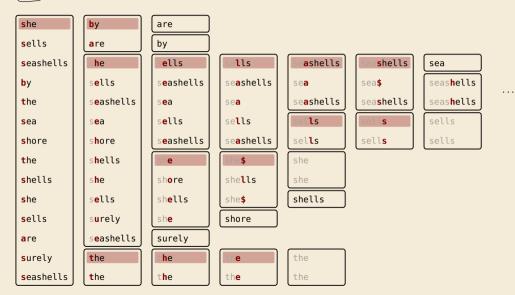
she	b y	are		
s ells	a re	by		
s eashells	she	sells	sells	seashells
b y	s e lls	s e ashells	se a shells	sea
the	s e ashells	s ea	se a	se a shells
sea	s ea	s e lls	se l ls	sel ls
shore	s h ore	s e ashells	se a shells	sel ls
the	s h ells	she	she\$	she
s hells	s he	shore	she lls	she
s he	s e lls	shells	she \$	shells
s ells	s <mark>u</mark> rely	sh e	shore	
a re	s e ashells	surely		
surely	the	the	the	the
seashells	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
 - → 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



- \rightarrow can show: $\sim 2 \ln(2) \cdot n \lg n \approx 1.39 n \lg n$ character comparisons on average
- simple to code
- ficient 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

random string

- **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
 - → never compare equal prefixes twice

for random strings

 \rightarrow can show: $\sim 2 \ln(2) \cdot n \lg n$ ≈ 1.39 $n \lg n$ character comparisons on average

- simple to code
- efficient for sorting many lists of strings
 - fat-pivot radix quicksort finds suffix array in $O(n \log n)$ expected time

for so the sorting $\Theta(L^2)$

T= a .- a\$

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