



ALGORITHMS OF BIOINFORMATICS

6

Suffix Trees

11 December 2025

Prof. Dr. Sebastian Wild

Outline

6 Suffix Trees

- 6.1 Suffix Trees
- 6.2 Direct Applications
- 6.3 Generalized Suffix Trees & Augmentation
- 6.4 Tandem Repeats
- 6.5 Longest Common Extensions
- 6.6 Suffix Arrays
- 6.7 Suffix Sorting: Induced Sorting and Merging
- 6.8 Suffix Sorting: The DC3 Algorithm
- 6.9 The LCP Array
- 6.10 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.

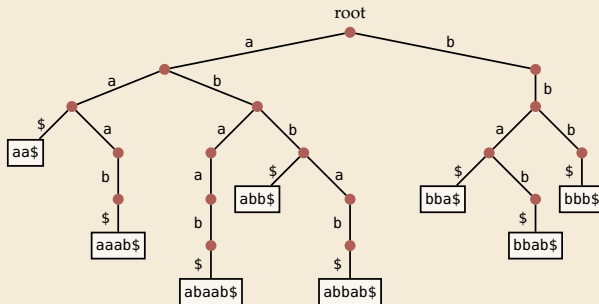
6.1 Suffix Trees

Recap: Tries

- ▶ efficient dictionary data structure for strings (or for Aho-Corasick automata 😊)
- ▶ name from **re**trieval, 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 \$ ✓

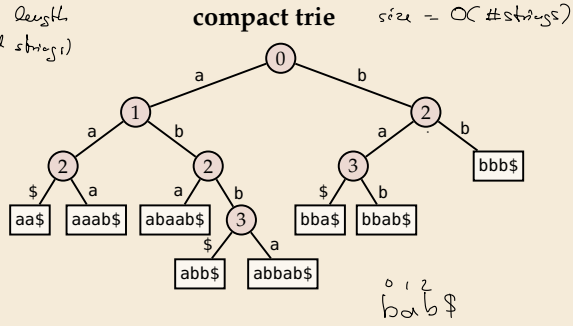
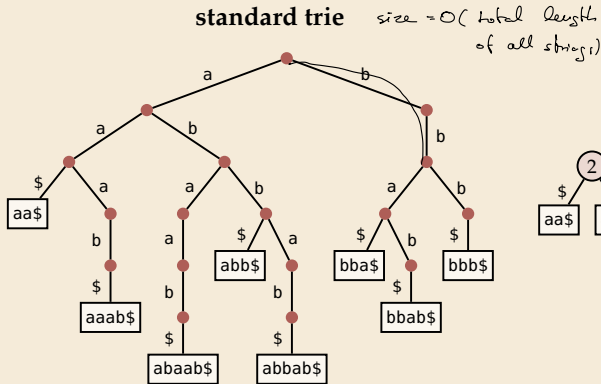
▶ **Example:**

{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



- ▶ search gives first character of edge only \rightsquigarrow must check for match against stored string
- ▶ all nodes ≥ 2 children \rightsquigarrow $\#nodes \leq \#leaves = \#strings \rightsquigarrow$ 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 $T[n] := \$$)

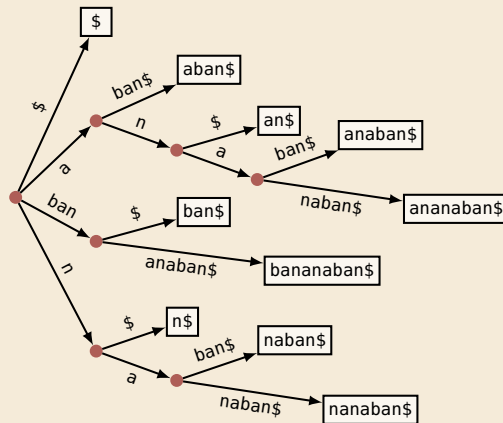
Example:

$T = \text{bananaban\$}$

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	\$

$T =$



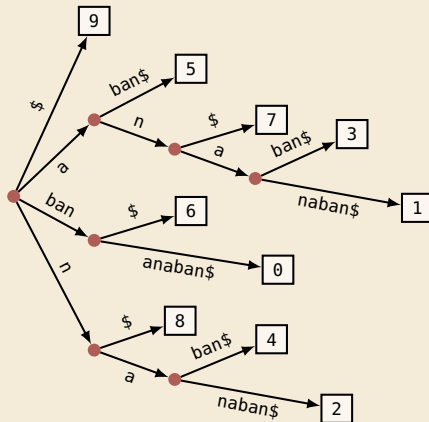
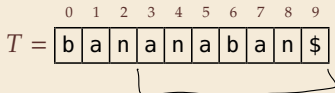
Suffix trees – Definition

- ▶ 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)

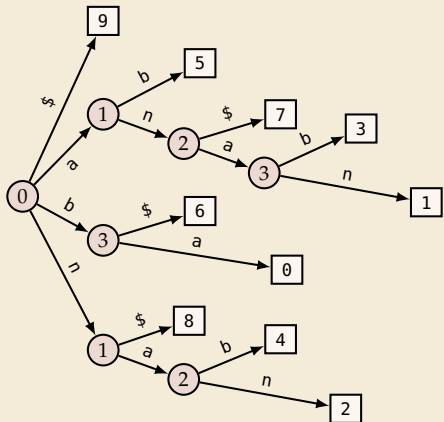
Example:

$T = \text{bananaban\$}$

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

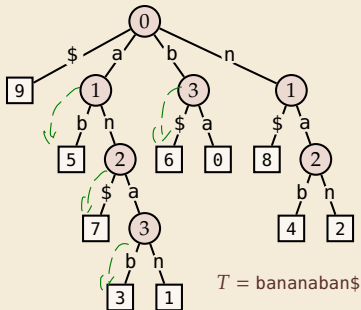
\rightsquigarrow for now, take linear-time construction for granted. What can we do with them?

6.2 Direct Applications

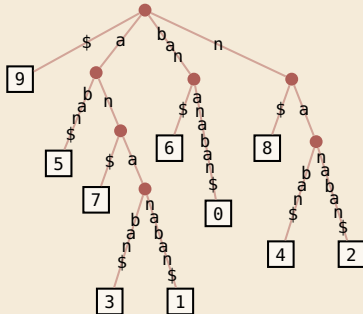
Applications of suffix trees

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

Recall: \mathcal{T} stored like this:



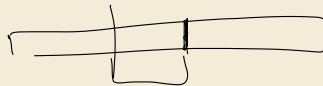
but think about this:



- 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

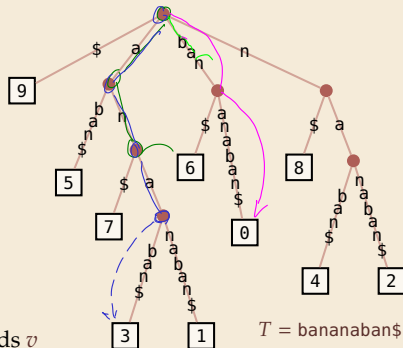
3. **we run out of tree**

reach a leaf 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:

► $P = \underline{\text{ann}}$

► $P = \underline{\text{baa}}$

► $P = \underline{\text{ana}}$

► $P = \text{ba}$

► $P = \underline{\text{briar}}$

Application 2: Longest repeated substring

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



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 $j \neq i$: $T[j..j + \ell) = T[i..i + \ell)$.



How can we efficiently check *all possible substrings*?



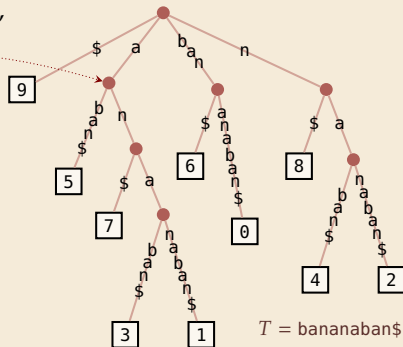
Repeated substrings = shared paths in *suffix tree*



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



Repeated substrings = shared paths in *suffix tree*



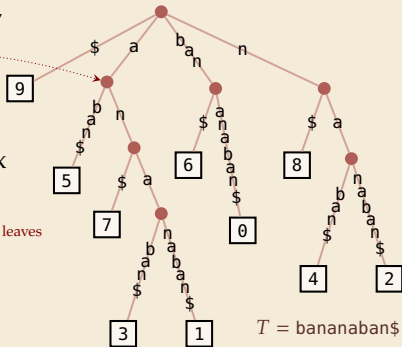
- $T_5 = \text{aban\$}$ and $T_7 = \text{an\$}$ have *longest common prefix* 'a'

$\rightsquigarrow \exists$ 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



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*?



Repeated substrings = shared paths in *suffix tree*



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

↪ \exists internal node with path label 'a'

here single edge, can be longer path

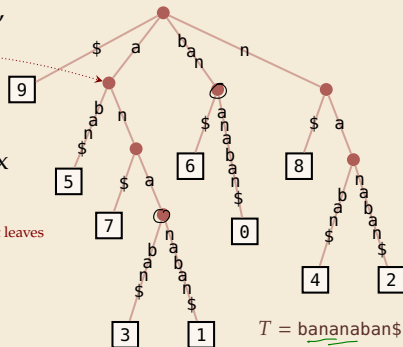
↪ longest repeated substring = longest common prefix (LCP) of two suffixes

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 ↪ $\Theta(n)$ time



6.3 Generalized Suffix Trees & Augmentation

Generalized suffix trees

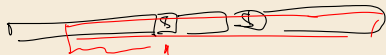
- ▶ longest *repeated* substring (of one string) feels very similar to longest *common* substring of several strings $T^{(1)}, \dots, T^{(k)}$ with $T^{(j)} \in \Sigma^{n_j}$
- ▶ can we solve that in the same way?
- ▶ could build the suffix tree for each $T^{(j)}$... but doesn't seem to help

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
- \rightsquigarrow 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*

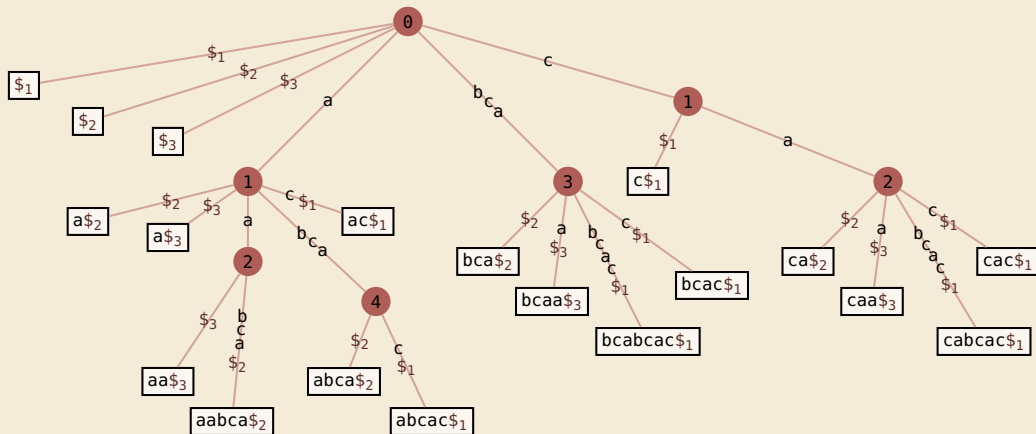
- ▶ Define $T := T^{(1)}\underline{\$}_1 T^{(2)}\underline{\$}_2 \dots T^{(k)}\underline{\$}_k$ for k new end-of-word symbols
- ▶ Construct suffix tree \mathcal{T} for T

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



Generalized Suffix Tree – Example

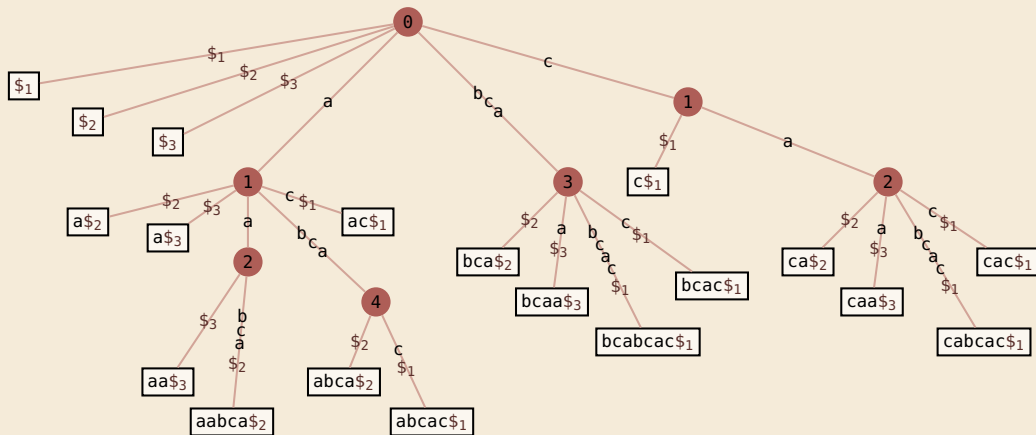
$T^{(1)} = \text{bcabcac}$, $T^{(2)} = \text{aabca}$, $T^{(3)} = \text{bcaa}$



Application 3: Longest common substring

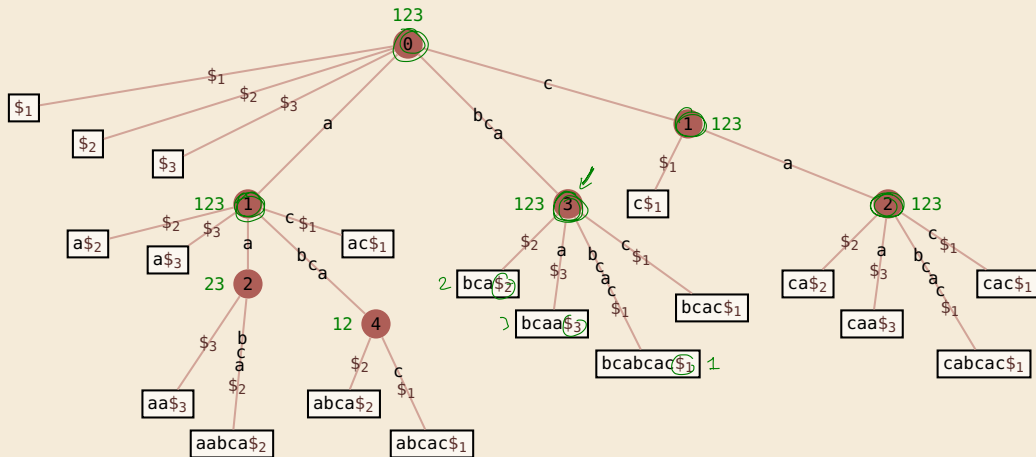
- ▶ With that new idea, we can find longest common substrings:
 1. Compute generalized suffix tree \mathcal{T} .
 2. Store with each node the *subset of strings* that contain its path label:
 - 2.1. Traverse \mathcal{T} bottom-up.
 - 2.2. For a leaf (j, i) , the subset is $\{j\}$.
 - 2.3. For an internal 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 $\{1, \dots, k\}$.
- ▶ Steps 1, 3, 4 take time $\Theta(n)$ for $n = n_1 + \dots + n_k$ the total length of all texts.
- ▶ Step 2 takes $\Theta(kn)$ time \rightsquigarrow linear if $k = O(1)$.
 - ▶ dependence on k can be removed with more clever labeling of vertices

Longest common substring – Example

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


Longest common substring – Example

$T^{(1)} = \text{bcab}\underline{\text{c}}\text{ac}$, $T^{(2)} = \text{aab}\underline{\text{c}}\text{a}$, $T^{(3)} = \text{bca}\underline{\text{a}}$



Application 4: DNA Sample Contamination

- ▶ **Given:** new reads $P[0..p)$, known contamination sources $C[0..c)$, length ℓ
- ▶ **Goal:** indices $I \subseteq [0..p)$ of reads that
do **not** have common substring of length $\geq \ell$
with **any** contamination source $C[j]$.

Application 4: DNA Sample Contamination

► **Given:** new reads $P[0..p)$, known contamination sources $C[0..c)$, length ℓ

► **Goal:** indices $I \subseteq [0..p)$ of reads that
do **not** have common substring of length $\geq \ell$
with **any** contamination source $C[j]$.

► Solution similar to longest common substring

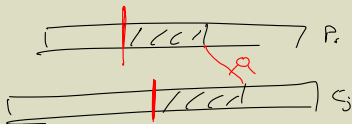
► Build generalized suffix tree \mathcal{T} of $P[0..p)$ and $C[0..c)$

► $P[r]$ has common substring with $C[j]$ of length $\geq \ell$
iff \exists node v in \mathcal{T} of string depth $\geq \ell$ and leaves from $P[r]$ and $C[j]$ below.

\rightsquigarrow Mark each node in \mathcal{T} in a bottom-up traversal if its subtree contains any leaves from any $C[j]$.

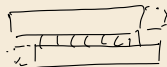
► For all marked nodes v of string depth $\geq \ell$, add r to I for all $P[r]$ with a leaf in subtree of v

\rightsquigarrow Linear in the total length of P and C .



Application 5: Computing the Overlap Graph

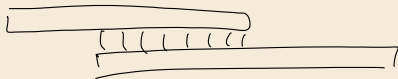
- Recall the genome sequencing problem with inexact reads



↪ cannot assume perfect coverage / $k - 1$ char overlaps for k -mers

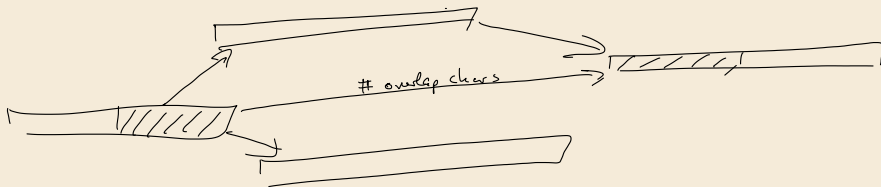
- Here: Overlap graph with edge weights = **length of exact suffix-prefix match**

↪ **All Pairs Suffix-Prefix Matching**



- **Given:** reads $P[0..p)$, $P[r] = P_r[0..m_r)$

- **Goal:** $O[0..p)[0..p)$ where $O[i][j]$ is length of longest suffix of P_i that is a prefix of P_j .



Application 5: Computing the Overlap Graph

- ▶ Recall the genome sequencing problem with inexact reads

↪ cannot assume perfect coverage / $k - 1$ char overlaps for k -mers

- ▶ Here: Overlap graph with edge weights = **length of exact suffix-prefix match**

↪ **All Pairs Suffix-Prefix Matching**

- ▶ **Given:** reads $P[0..p)$, $P[r] = P_r[0..m_r)$

- ▶ **Goal:** $O[0..p)[0..p)$ where $O[i][j]$ is length of longest suffix of P_i that is a prefix of P_j .

- ▶ To find suffix-prefix overlaps, use generalized suffix tree \mathcal{T}

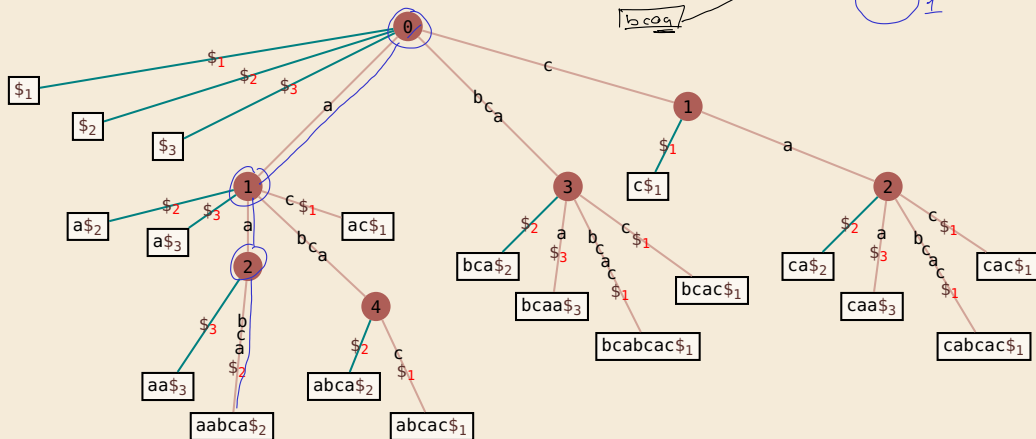
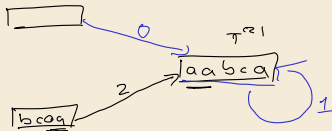
- ▶ Mark edges labeled only with $\$_i$ as *terminal edges*

- ▶ Suffix of P_i equals prefix of P_j iff traversal of P_j in \mathcal{T} reaches $\$_j$ terminal edge

↪ in a single traversal, remember deepest terminal edge for each string
output all when reaching leaf for $P_j[0..m_j)$

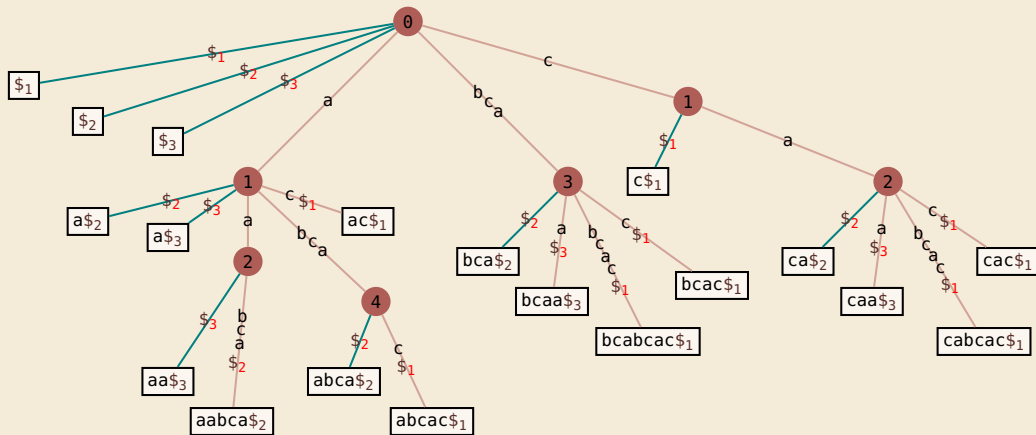
All-Pairs Suffix-Prefix – Example

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



terminal edge

All-Pairs Suffix-Prefix – Example

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


terminal edge

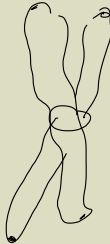
6.4 Tandem Repeats

Repetitions and “Junk DNA”

- ▶ only 1–2% of human genome are protein-coding genes . . . *what is the rest there for?*

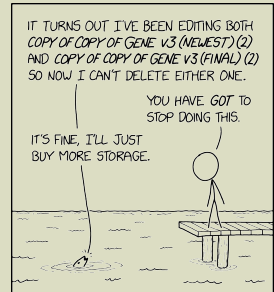
Repetitions and “Junk DNA”

- ▶ only 1–2% of human genome are protein-coding genes . . . *what is the rest there for?*
- ▶ some 15-25% are actually functional in that they are involved in some direct actions
 - ▶ protein-coding genes
 - ▶ regions representing non-coding RNA (*transfer RNA, ribosomal RNA, interfering RNA, . . .*)
 - ▶ structural support for chromosomes (*telomeres, centromeres, satellite DNA*)



Repetitions and “Junk DNA”

- ▶ only 1–2% of human genome are protein-coding genes . . . *what is the rest there for?*
- ▶ some 15-25% are actually functional in that they are involved in some direct actions
 - ▶ protein-coding genes
 - ▶ regions representing non-coding RNA (*transfer RNA, ribosomal RNA, interfering RNA, . . .*)
 - ▶ structural support for chromosomes (*telomeres, centromeres, satellite DNA*)
- ▶ *how about the rest?*
 - ▶ a lot is *transposons* (“jumping genes”, “mobile DNA”) function is unclear each copy leaves a few newly inserted bases in your genome . . .
 - ▶ *introns* (cut out parts in eukaryotic split genes)
 - ▶ *pseudogenes* former copies of genes that lost function due to mutation

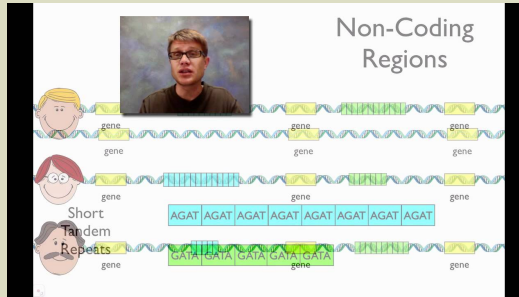


WHY LUNGFISH HAVE SUCH ENORMOUS GENOMES

xkcd.com/3864/

DNA Fingerprinting

- ▶ for identifying individuals, need highly variable DNA region
- ▶ but also one that we can easily find in reads



▶ DNA Fingerprinting
<https://youtu.be/DbR9xMXuK7c>

*How can we find **tandem repeats**?*

Tandem Repeats

Formally, finding tandem repeats amounts to the following problem.

× ×

- ▶ **Given:** Text $T[0..n)$
- ▶ **Goal:** Report all pairs (i, ℓ) , so that $T[i..i + \ell) = T[i + \ell, i + 2\ell)$

Tandem Repeats

Formally, finding tandem repeats amounts to the following problem.

- ▶ **Given:** Text $T[0..n)$
- ▶ **Goal:** Report all pairs (i, ℓ) , so that $T[i..i + \ell) = T[i + \ell, i + 2\ell)$

How can suffix trees help us here?

Tandem Repeats

Formally, finding tandem repeats amounts to the following problem.

- ▶ **Given:** Text $T[0..n)$
- ▶ **Goal:** Report all pairs (i, ℓ) , so that $T[i..i + \ell) = T[i + \ell, i + 2\ell)$

How can suffix trees help us here? They can't (directly) 🙄

But remember the *longest common extension (LCE)* data structure:

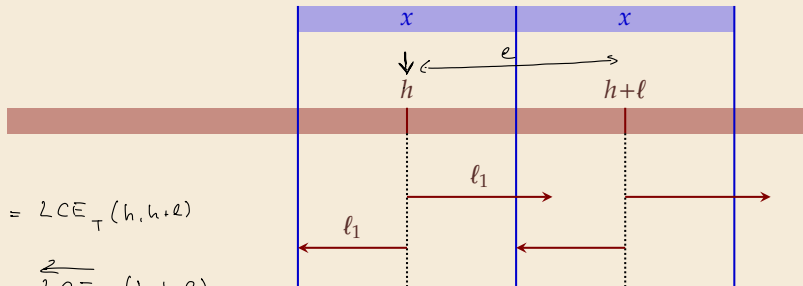
- ▶ **Given:** String $T[0..n)$
- ▶ **Goal:** Answer queries $\text{LCE}_T(i, j) = \max\{\ell : T[i..i + \ell) = T[j..j + \ell)\}$



After $\Theta(n)$ preprocessing (time and space), query time $O(1)$

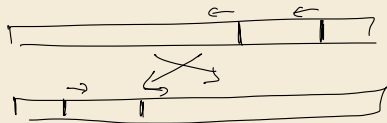
Tandem Repeats – Idea

- ▶ a tandem repeat is a substring xx
- ▶ if we fix the length $\ell = |x|$ and a “seed position” h ,
a tandem repeat must allow extending outwards from h by at least ℓ positions



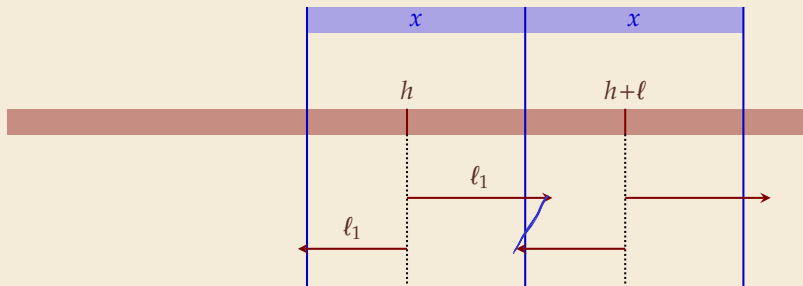
$$\ell_1 = \text{LCE}_T(h, h+\ell)$$

$$\ell_2 = \text{LCE}_T(h, h+\ell) = \text{LCE}_{TR}(n-h, n-h-\ell)$$

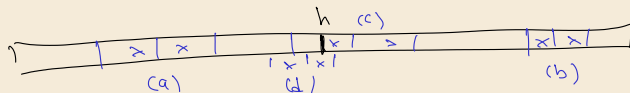


Tandem Repeats – Idea

- ▶ a tandem repeat is a substring xx
- ▶ if we fix the length $\ell = |x|$ and a “seed position” h ,
a tandem repeat must allow extending outwards from h by at least ℓ positions

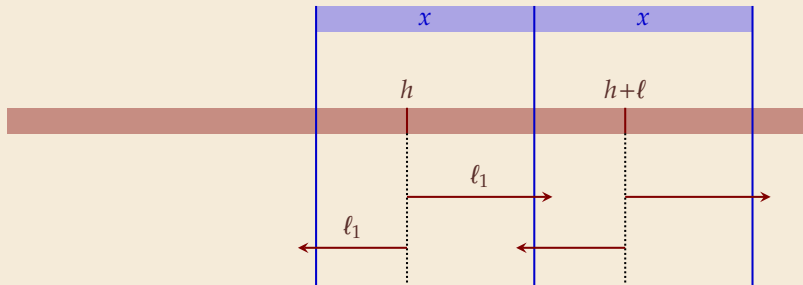


D & C!



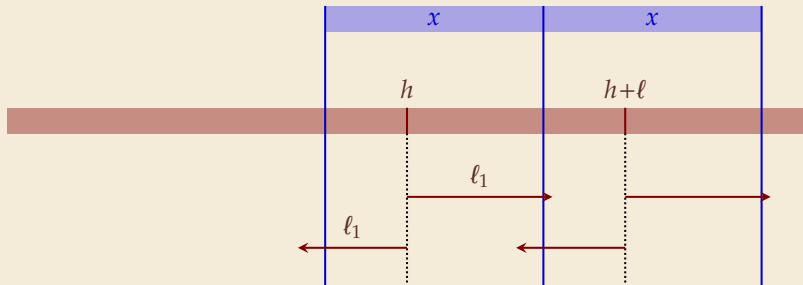
Tandem Repeats – Idea

- ▶ a tandem repeat is a substring xx
- ▶ if we fix the length $\ell = |x|$ and a “seed position” h ,
a tandem repeat must allow extending outwards from h by at least ℓ positions



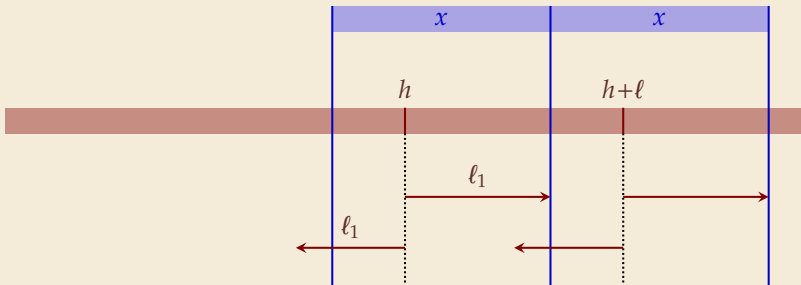
Tandem Repeats – Idea

- ▶ a tandem repeat is a substring xx
- ▶ if we fix the length $\ell = |x|$ and a “seed position” h ,
a tandem repeat must allow extending outwards from h by at least ℓ positions



Tandem Repeats – Idea

- ▶ a tandem repeat is a substring xx
- ▶ if we fix the length $\ell = |x|$ and a “seed position” h ,
a tandem repeat must allow extending outwards from h by at least ℓ positions



Tandom Repeats – Code

```
1 procedure tandemRepeats( $T[0..n]$ )
2   Build LCE data structure on  $T$  and  $T^R$ 
3   //  $T^R = \text{reversed text}$ 
4   // We abbreviate  $\vec{e}(i, j) = \text{LCE}_T(i, j)$ 
5   // resp.  $\tilde{e}(i, j) = \text{LCE}_{T^R}(n - j, n - i)$ 
6    $tr := \emptyset$ 
7   appendRepeats( $T[0..n]$ ,  $tr$ )
8   return  $tr$ 
```

```
1 procedure appendRepeats( $T[0..n]$ ,  $tr$ )
2   if  $n \leq 1$  return
3    $h := \lfloor n/2 \rfloor$ 
4   appendRepeats( $T[0..h]$ )
5   appendRepeats( $T[h..n]$ )
6   // Add all repeats  $xx$  with  $h$  inside first copy of  $x$ 
7   for  $\ell \in [1 .. n - h]$  // length of  $x$ 
8      $q := h + \ell$ 
9      $\ell_1 := \vec{e}(h, q)$ 
10     $\ell_2 := \tilde{e}(h - 1, q - 1)$ 
11    if  $\ell_1 + \ell_2 \geq \ell$ 
12       $tr := tr \cup \{(j, \ell) : j \in [h - \ell_2 .. h + \ell_1 - \ell]\}$ 
13    // Symmetrically for  $h$  inside second copy of  $x$ 
14    for  $\ell \in [1..h]$ 
15       $q := h - \ell$ 
16       $\ell_1 := \vec{e}(h, q)$ 
17       $\ell_2 := \tilde{e}(h - 1, q - 1)$ 
18      if  $\ell_1 + \ell_2 \geq \ell$ 
19         $tr := tr \cup \{(j, \ell) : j \in [q - \ell_2 .. q + \ell_1 - \ell]\}$ 
```

Tandom Repeats – Analysis

► Running time

- nonrecursive part $O(n)$ (using LCE preprocessing)
- recording output pairs (j, ℓ) over all recursive calls $O(\text{output})$ time
- recursive function uses $O(n)$ LCE queries $\rightsquigarrow O(n)$ time
- recursive function satisfies $T(n) = \Theta(n) + 2T(n/2) \rightsquigarrow T(n) = O(n \log n)$

$\rightsquigarrow O(n \log n + \text{output})$

6.5 Longest Common Extensions

Longest Common Extensions & Suffix Trees

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

Recall *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: $\text{LCE}(i, j) = \max\{\ell : T[i..i + \ell) = T[j..j + \ell)\}$

Longest Common Extensions & Suffix Trees

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

Recall *longest common extension (LCE)* data structure

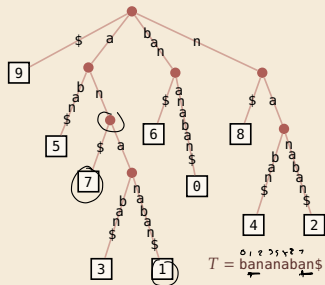
- ▶ **Given:** String $T[0..n)$
- ▶ **Goal:** Answer LCE queries, i. e.,
 given positions i, j in T ,
 how far can we read the same text from there?
 formally: $\text{LCE}(i, j) = \max\{\ell : T[i..i + \ell) = T[j..j + \ell)\}$

↪ use suffix tree of T !

(length of) longest common prefix
of i th and j th suffix

- In \mathcal{T} : $\text{LCE}(i, j) = \text{LCP}(T_i, T_j) \rightsquigarrow$ same thing, different name!
 $=$ string depth of
lowest common ancestor (LCA) of
leaves \boxed{i} and \boxed{j}

- in short: $\text{LCE}(i, j) = \text{LCP}(T_i, T_j) = \text{stringDepth}(\text{LCA}(\boxed{i}, \boxed{j}))$



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 🙄

Efficient LCA

How to find lowest common ancestors?

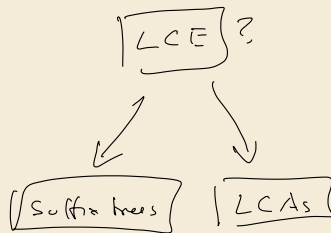
- ▶ Could walk up the tree to find LCA $\rightsquigarrow \Theta(n)$ worst case 🙄
- ▶ Could store all LCAs in big table $\rightsquigarrow \Theta(n^2)$ space and preprocessing 🙄



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

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

and suffix tree construction inside ...



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

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

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



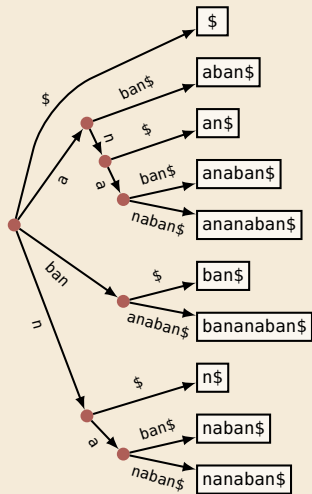
many pointers in tree incur large space overhead



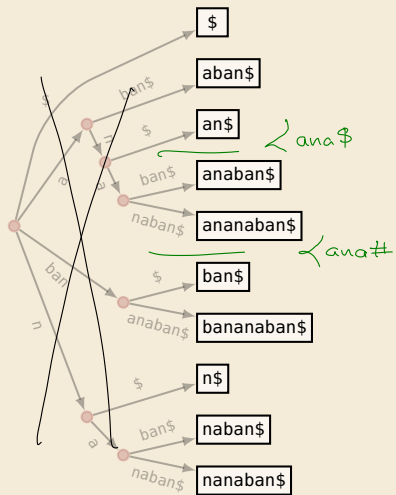
6.6 Suffix Arrays

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*

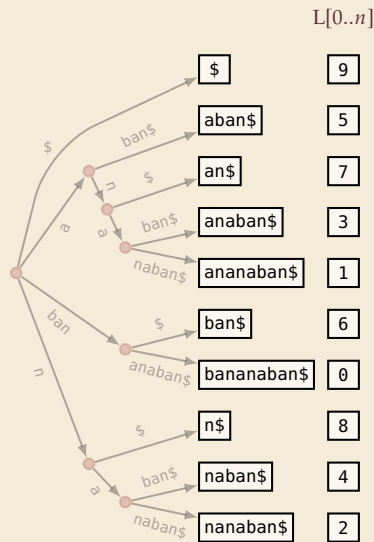
suffix array
//

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

ana\$
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
 $\leq (\lg n + 2) \cdot \underline{m}$ character comparisons

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

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!

Excursion: String sorting

- ▶ when sorting strings, “blind” ^{string} comparisons can cost $\Theta(n)$ character comparisons
- ▶ happens iff strings share long prefix!
- ↪ dedicated string sorting methods need to remember common prefixes between strings
then we can avoid redoing these comparisons

Excursion: String sorting

- ▶ when sorting strings, “blind” comparisons can cost $\Theta(n)$ character comparisons
- ▶ happens iff strings share long prefix!
- ↪ dedicated string sorting methods need to remember common prefixes between strings then we can avoid redoing these comparisons

(length of) longest common prefix

- ▶ Option 1: Mergesort with LCP values for adjacent elements in runs *↪ exam*
- ▶ Option 2: Fat-pivot radix quicksort
 - ▶ **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

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

Fat-pivot radix quicksort – Example

she

sells

seashells

by

the

sea

shore

the

shells

she

sells

are

surely

seashells

Fat-pivot radix quicksort – Example

she
sells
seashells
by
the
sea
shore
the
shells
she
sells
are
surely
seashells

Fat-pivot radix quicksort – Example

she	by
sells	are
seashells	she
by	sells
the	seashells
sea	sea
shore	shore
the	shells
shells	she
she	sells
sells	surely
are	seashells
surely	the
seashells	the

Fat-pivot radix quicksort – Example

she	by	are
sells	are	by
seashells	she	
by	sells	
the	seashells	
sea	sea	
shore	shore	
the	shells	
shells	she	
she	sells	
sells	surely	
are	seashells	
surely	the	
seashells	the	

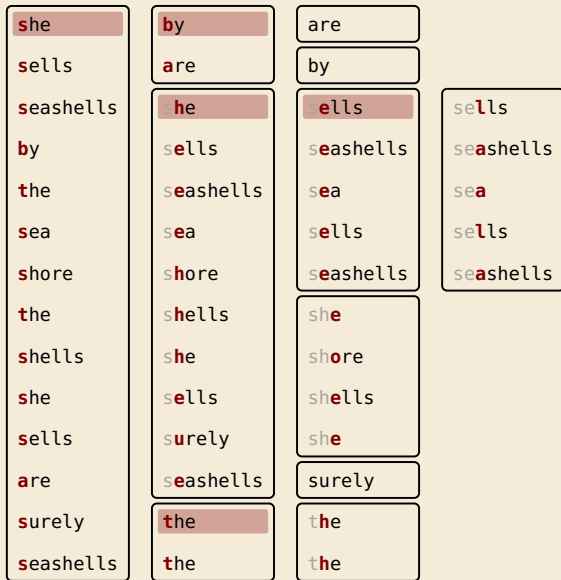
Fat-pivot radix quicksort – Example

she	by	are
sells	are	by
seashells	he	sells
by	sells	seashells
the	seashells	sea
sea	sea	sells
shore	shore	seashells
the	shells	she
shells	she	shore
she	sells	shells
sells	surely	she
are	seashells	surely
surely	the	
seashells	the	

Fat-pivot radix quicksort – Example

she	by	are
sells	are	by
seashells	he	sells
by	sells	seashells
the	seashells	sea
sea	sea	sells
shore	shore	seashells
the	shells	she
shells	she	shore
she	sells	shells
sells	surely	she
are	seashells	surely
surely	the	the
seashells	the	the

Fat-pivot radix quicksort – Example



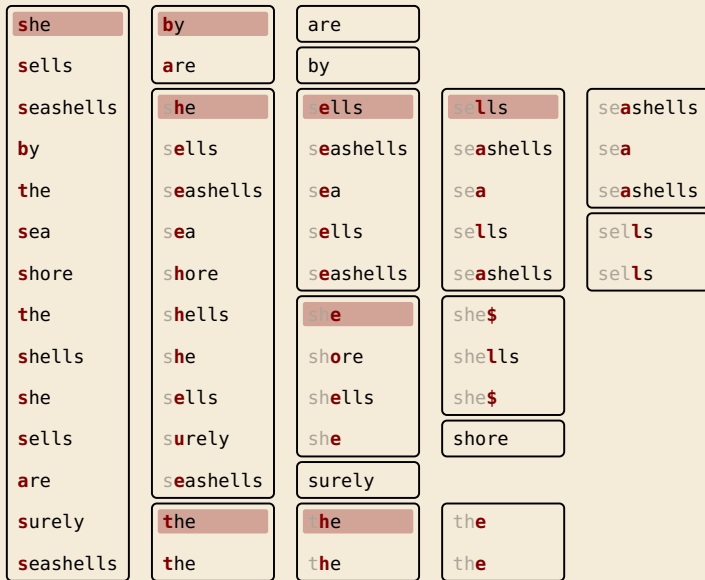
Fat-pivot radix quicksort – Example

she	by	are	
sells	are	by	
seashells	he	sells	sells
by	sells	seashells	seashells
the	seashells	sea	sea
sea	sea	sells	sells
shore	shore	seashells	seashells
the	shells	she	she\$
shells	she	shore	shells
she	sells	shells	she\$
sells	surely	she	shore
are	seashells	surely	
surely	the	the	
seashells	the	the	

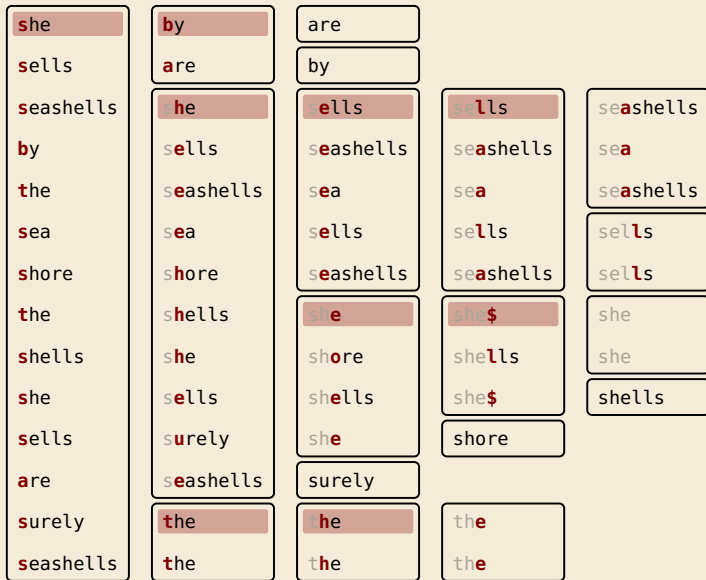
Fat-pivot radix quicksort – Example

she	by	are	
sells	are	by	
seashells	he	sells	sells
by	sells	seashells	seashells
the	seashells	sea	sea
sea	sea	sells	sells
shore	shore	seashells	seashells
the	shells	she	she\$
shells	she	shore	shells
she	sells	shells	she\$
sells	surely	she	shore
are	seashells	surely	
surely	the	the	the
seashells	the	the	the

Fat-pivot radix quicksort – Example



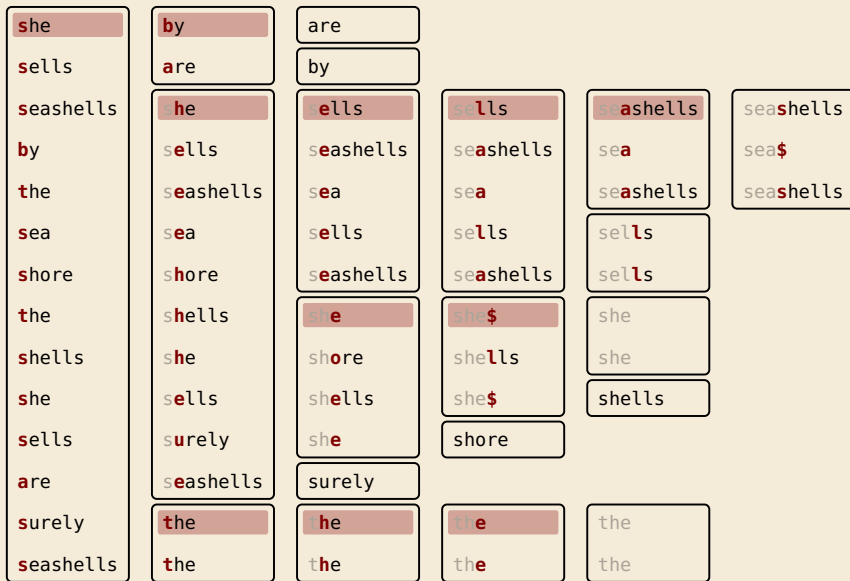
Fat-pivot radix quicksort – Example



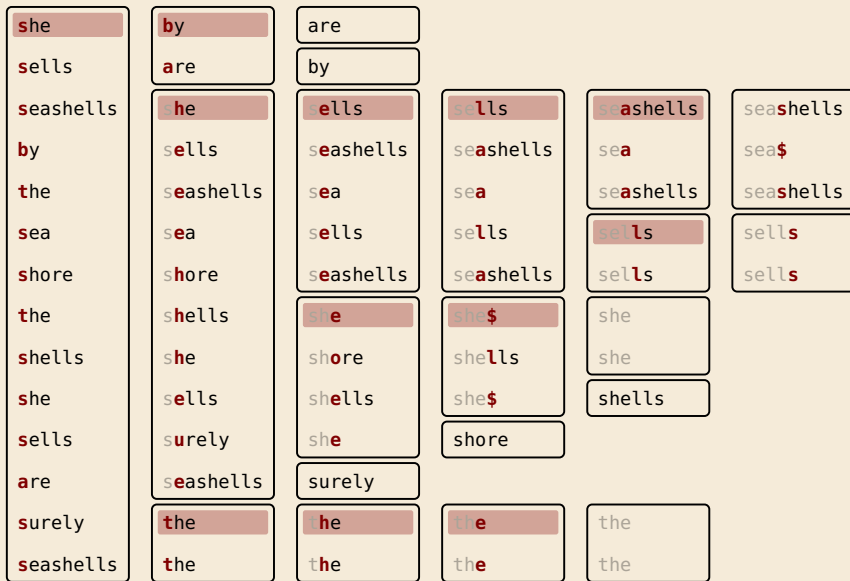
Fat-pivot radix quicksort – Example

she	by	are		
sells	are	by		
seashells	he	sells	sells	seashells
by	sells	seashells	seashells	sea
the	seashells	sea	sea	seashells
sea	sea	sells	sells	sells
shore	shore	seashells	seashells	sells
the	shells	she	she\$	she
shells	she	shore	shells	she
she	sells	shells	she\$	shells
sells	surely	she	shore	
are	seashells	surely		
surely	the	the	the	the
seashells	the	the	the	the

Fat-pivot radix quicksort – Example



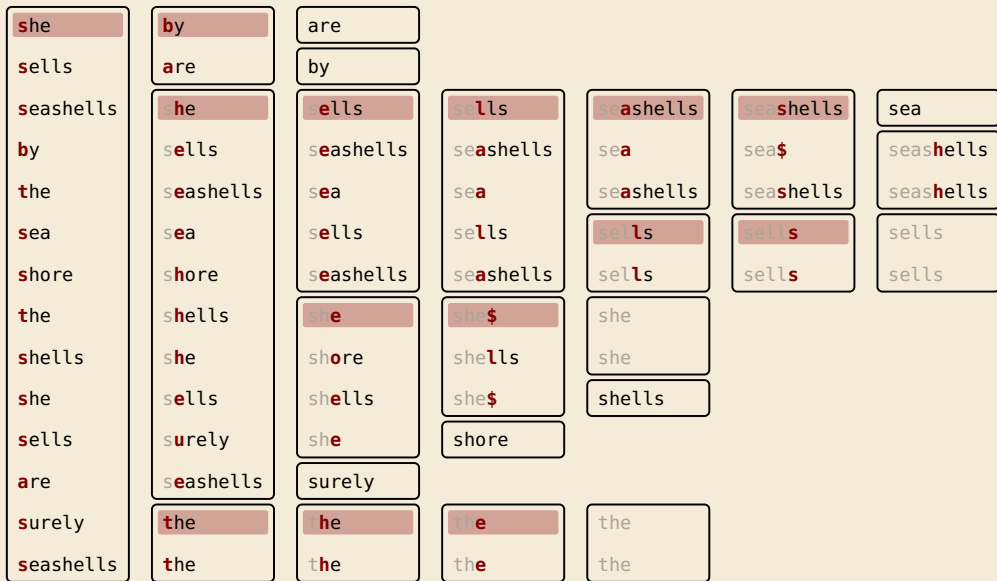
Fat-pivot radix quicksort – Example



Fat-pivot radix quicksort – Example

she	by	are				
sells	are	by				
seashells	she	sells	sells	seashells	seashells	sea
by	sells	seashells	seashells	sea	sea\$	seashells
the	seashells	sea	sea	seashells	seashells	seashells
sea	sea	sells	sells	sells	sells	
shore	shore	seashells	seashells	sells	sells	
the	shells	she	she\$	she	she	
shells	she	shore	shells	shells	shells	
she	sells	shells	she\$			
sells	surely	she	shore			
are	seashells	surely				
surely	the	the	the	the		
seashells	the	the	the	the		

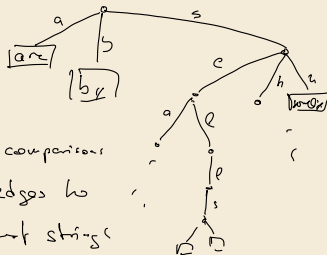
Fat-pivot radix quicksort – Example



Fat-pivot radix quicksort – Example

she	by	are					
sells	are	by					
seashells	she	sells	sells	seashells	seashells	seashells	sea
by	sells	seashells	seashells	sea	sea	seashells	seashells
the	seashells	sea	sea	sells	seashells	seashells	seashells
sea	sea	sells	sells	sells	sells	sells	sells
shore	shore	seashells	seashells	sells	sells	sells	sells
the	shells	she	she\$	she	she	she	
shells	she	shore	shells	shells	shells	shells	
she	sells	shells	she\$				
sells	surely	she	shore				
are	seashells	surely					
surely	the	the	the	the	the	the	
seashells	the	the	the	the	the	the	

#LCP comparisons
 = # edges to insert strings



Fat-pivot radix quicksort – Analysis

Separately analyze character comparisons **by outcome**

1. *“Decisive Comparisons:”* character comparisons with outcome “<” or “>”

▶ can have at most one in any **string comparison** (afterwards done!)

↪ Same number of decisive comparisons as in standard quicksort (just delayed)

↪ expected $\sim 2 \ln(2) \cdot n \lg n \approx 1.39n \lg n$ decisive comparisons

Fat-pivot radix quicksort – Analysis

Separately analyze character comparisons **by outcome**

1. *“Decisive Comparisons:”* character comparisons with outcome “<” or “>”

- ▶ can have at most one in any **string comparison** (afterwards done!)
- ↪ Same number of decisive comparisons as in standard quicksort (just delayed)
- ↪ expected $\sim 2 \ln(2) \cdot n \lg n \approx 1.39n \lg n$ decisive comparisons

2. *LCP comparisons:* character comparisons that return “=”

- ▶ must be all remaining character comparisons
- ▶ every such comparison contributes to common prefix between strings, never compare same characters again
- ▶ every string sort must discover longest common prefixes in sorted order
- ↪ #LCP comparisons = #comparisons when inserting all strings into a **trie**

(
w/ unused paths
↳ leaves collapsed

Fat-pivot radix quicksort – Discussion

👍 simple to code

👍 efficient for sorting many lists of strings

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

random string



👎 worst case remains $\Omega(n^2)$, i. e., $T = a^n$

Note: Not quicksort's fault! Any generic string sorting method must take $\Omega(n^2)$ time here

Fat-pivot radix quicksort – Discussion

👍 simple to code

👍 efficient for sorting many lists of strings

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

random string



👎 worst case remains $\Omega(n^2)$, i. e., $T = a^n$

Note: Not quicksort's fault! Any generic string sorting method must take $\Omega(n^2)$ time here

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