



6 Text Indexing – Searching whole genomes

16 March 2021

Sebastian Wild

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
- 6.7 The LCP Array

6.1 Motivation

Text indexing

- ▶ *Text indexing* (also: *offline text search*):

- ▶ case of string matching: find $P[0..m]$ in $T[0..n]$

- ▶ but with *fixed* text \rightsquigarrow preprocess T (instead of P)

- \rightsquigarrow expect many queries P , answer them without looking at all of T

- \rightsquigarrow essentially a data structuring problem: “building an *index* of T ”

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 \mapsto page numbers where they occur
 - ▶ assumption: searches are only for **whole (key) words**
- \rightsquigarrow often reasonable for natural language text

Inverted indices

same as "indexes"

- ▶ original indices in books: list of (key) words \mapsto page numbers where they occur
 - ▶ assumption: searches are only for **whole** (key) **words**
- \rightsquigarrow 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)
goes \rightarrow go, cats \rightarrow cat
- ▶ store mapping from words to a list of occurrences \rightsquigarrow how?

like a dictionary!

keys = words

values = list of occurrences,

BST/D
but $O(\log n)$
time

Clicker Question



Do you know what a *trie* is?

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

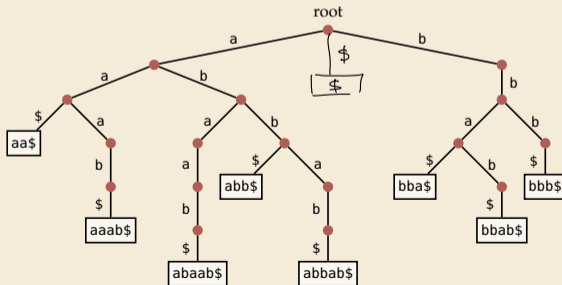
sli.do/comp526

Click on "Polls" tab

Tries

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

- ▶ **Example:** $\Sigma = \{a, b\}$
{aa\$, aaab\$, abaab\$, abb\$,
abbab\$, bba\$, bbab\$, bbb\$, \$}



Clicker Question

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

We now search for a query string Q with $|Q| = q$. ($q \leq m$)

How many **nodes** in the trie are **visited** during this **query**?



A $\Theta(\log n)$

F $\Theta(\log m)$

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

G $\Theta(q)$

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

H $\Theta(\log q)$

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

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

E $\Theta(m)$

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

sli.do/comp526

Click on "Polls" tab

Clicker Question

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

We now search for a query string Q with $|Q| = q$.

How many **nodes** in the trie are **visited** during this **query**?



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

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

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

G $\Theta(q)$ ✓

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

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

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

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

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

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

sli.do/comp526

Click on "Polls" tab

Clicker Question



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

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

A $\Theta(n)$

D $\Theta(n \log m)$

B $\Theta(n + m)$

E $\Theta(m)$

C $\Theta(n \cdot m)$

F $\Theta(m \log n)$

sli.do/comp526

Click on "Polls" tab

Clicker Question



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

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

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

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

B ~~$\Theta(n + m)$~~

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


C $\Theta(n \cdot m)$ ✓


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


sli.do/comp526

Click on "Polls" tab

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

```
ACAAGATGCCATTGTCCCCGGCCTCCTGCTGCTGCTGCTCTCCGGGGCCACGGCCACCGCTGCCCTGCCCTGGAGGGTGGCCCCACCGGC  
CGAGACAGCGAGCATATGCAGGAAGCGGCAGGAATAAGGAAAAGCAGCTCCTGACTTTCTCGCTTGGTGGTTTGTAGTGGACCTCCAGGC  
CAGTGCCGGGCCCTCATAGGAGAGGAAGCTCGGGAGGTGGCCAGGCGGCAGGAAGGCGCACCCCCAGCAATCCGCGCGCCGGGACAGAA  
TGCCCTGCAGGAACCTTCTTCTGGAAGACCTTCTCCTCTGCAAATAAAACCTCACCCATGAATGCTCACGCAAGTTTAATTACAGACCTGAA
```

- ▶ binary streams

```
00000010101001111010111000001111100011111011111001101101000011100010011011110000010001101010  
011011000011010110100000001000000011101011000001000011110101110110010001100101101110111111  
11000101000101100101000000111010101001100000001101100001100111110000101 010101110111100011  
1010111001001010101010000011111010011000000111100110101000000100100100000101100011000110111
```

~> need new ideas

6.2 Suffix Trees

Suffix trees – A ‘magic’ data structure

Appetizer: Longest common substring problem

- ▶ Given: strings S_1, \dots, S_k **Example:** $S_1 = \text{superiorcalifornialives}$, $S_2 = \text{sealiver}$
- ▶ 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 \rightsquigarrow 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 $T[n] := \$$)

0
banana \$

1
anana \$

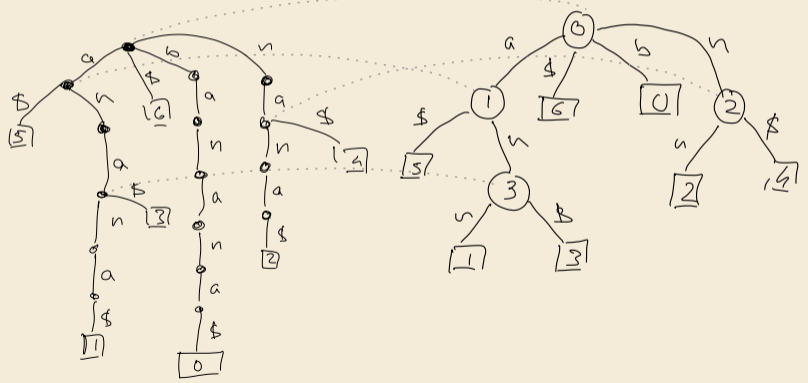
2
nana \$

3
ana \$

4
na \$

5
a \$

6
\$



Suffix trees – Definition

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

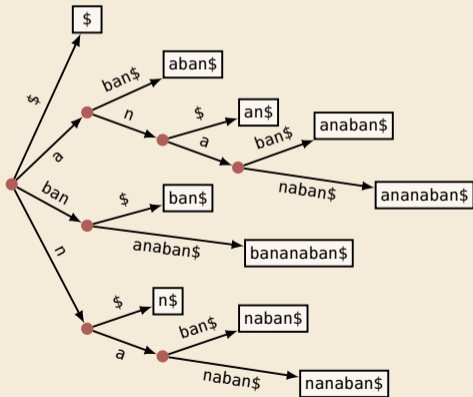
Example:

$T = \text{bananaban\$}$

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

$T =$

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



Suffix trees – Definition

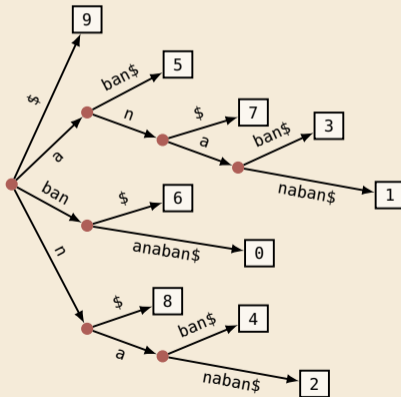
- ▶ suffix tree \mathcal{T} for text $T = T[0..n)$ = compact trie of all suffixes of $T\$$ (set $T[n] := \$$)
- ▶ except: in leaves, store *start index* (instead of actual string)

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	\$



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 actual string)

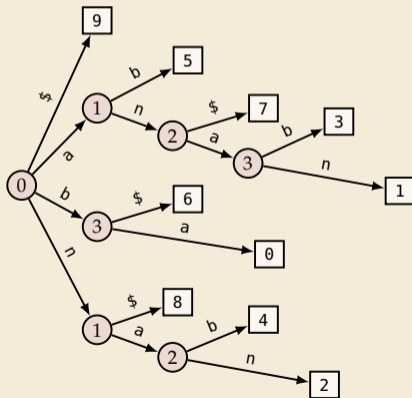
Example:

$T = \text{bananaban\$}$

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
 $T = \boxed{\text{b}} \boxed{\text{a}} \boxed{\text{n}} \boxed{\text{a}} \boxed{\text{n}} \boxed{\text{a}} \boxed{\text{b}} \boxed{\text{a}} \boxed{\text{n}} \boxed{\text{\$}}$

- ▶ 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 $\underbrace{\Theta(n^2)}$. \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!

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

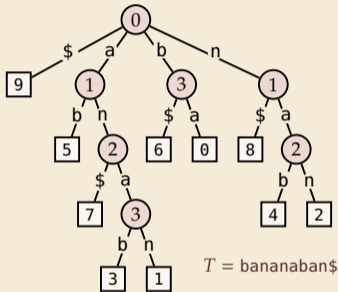
\rightsquigarrow for now, take linear-time construction for granted. What can we do with them?

6.3 Applications

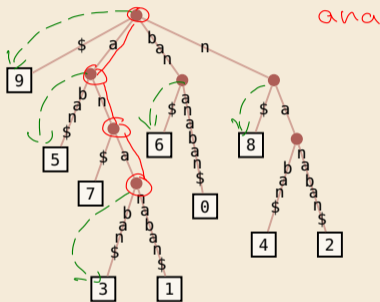
Applications of suffix trees

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

Recall: \mathcal{T} stored like this:

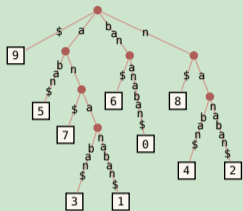


but think about this:



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

Clicker Question



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

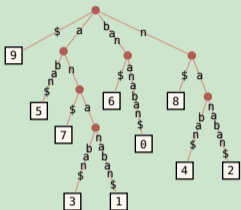
Check all that apply to this example.

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

sli.do/comp526

Click on "Polls" tab

Clicker Question



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.

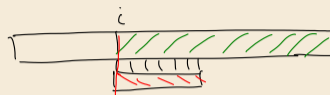
- A** ~~Nothing.~~
- B** P occurs in T . ✓
- C** ~~P does not occur in T .~~
- D** ~~P occurs once in T .~~
- E** P occurs twice in T . ✓
- F** ~~P starts at index 0.~~
- G** P starts at index 1. ✓
- H** ~~P starts at index 2.~~
- I** P starts at index 3. ✓
- J** ~~P starts at index 4.~~
- K** ~~P starts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

sli.do/comp526

Click on "Polls" tab

Application 1: Text Indexing / String Matching

- ▶ P occurs in $T \iff \underline{P}$ 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

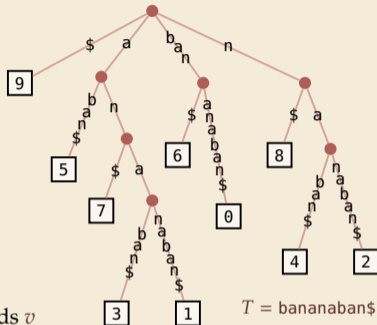
3. **we run out of tree**

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:

► $P = \text{ann}$

► $P = \text{ana}$

► $P = \text{briar}$

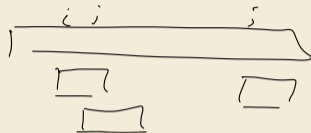
Application 2: Longest repeated substring

- **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*?

e. g. for compression \rightsquigarrow Unit 7



Application 2: Longest repeated substring

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

e.g. for compression \rightsquigarrow Unit 7



How can we efficiently check *all possible substrings*?



Repeated substrings = shared paths in *suffix tree*



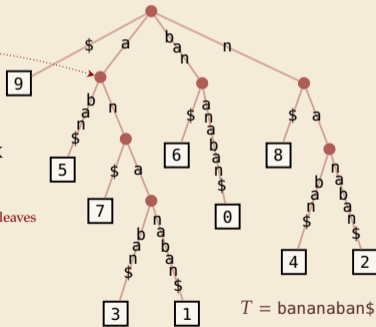
► $T_5 = \text{aban}\$$ and $T_7 = \text{an}\$$ have *longest common prefix* 'a'

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

here single edge, can be longer path

\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*?

e.g. for compression \rightsquigarrow Unit 7



Repeated substrings = shared paths in *suffix tree*



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

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

here single edge, can be longer path

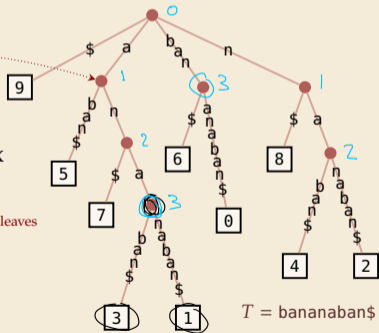
\rightsquigarrow longest repeated substring = longest common prefix (LCP) of two suffixes

actually: adjacent leaves

- ▶ **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 $T^{(1)}, \dots, T^{(k)}$ with $T^{(j)} \in \Sigma^{n_j}$
- ▶ can we solve that in the same way?
- ▶ could build the suffix tree for each $T^{(j)}$... but doesn't seem to help

Generalized suffix trees

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

Generalized suffix trees

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

Enter: *generalized suffix tree*

- ▶ Define $T := T^{(1)}\$1 T^{(2)}\$2 \dots T^{(k)}\$k$ 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 $T_i^{(j)} = T^{(j)}[i..n_j]$



Clicker Question



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

sli.do/comp526

Click on "Polls" tab

A large, light pink arrow with a white outline points to the right, starting from the text 'Click on "Polls" tab' and extending towards the right edge of the slide.

Application 3: Longest common substring

- ▶ With that new idea, we can find longest common superstrings:
 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 $\{1, \dots, k\}$.

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

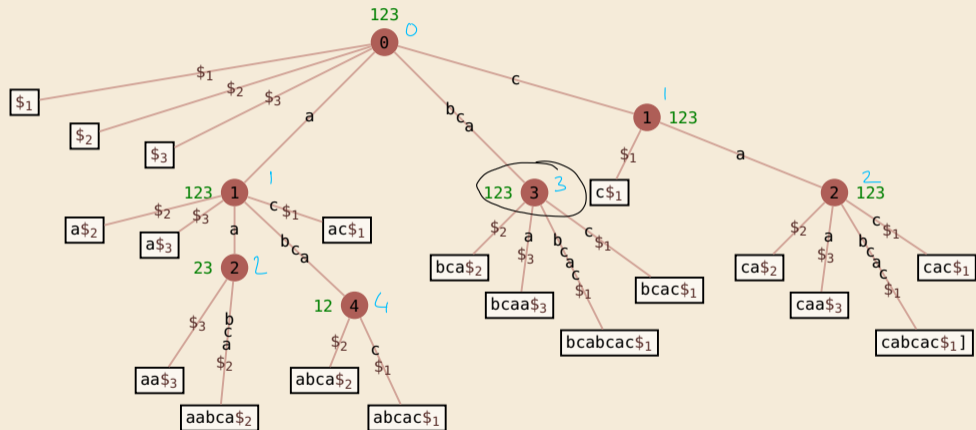
⊙ stores set of j so that there is a leaf (j, i) in the subtree

“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

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



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:

- ▶ **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: $LCE(i, j) = \max\{\ell : T[i..i + \ell] = T[j..j + \ell]\}$

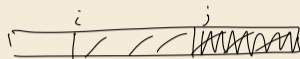


Application 4: Longest Common Extensions

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

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

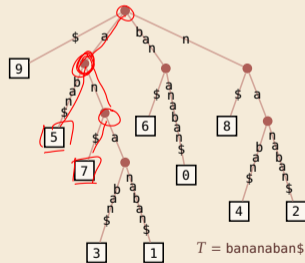
- ▶ **Given:** String $T[0..n]$
- ▶ **Goal:** Answer LCE queries, i. e.,
 given positions i, j in T ,
 how far can we read the same text from there?
 formally: $LCE(i, j) = \max\{\ell : T[i..i + \ell] = T[j..j + \ell]\}$



$\underline{a}ban\$$
 $\quad \underline{a}n\$$

↪ use suffix tree of T !

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





$T = \text{bananaban}\$$

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



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 LCA 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:**

- ▶ smallest i so that $T[i..i+m)$ and P differ in at most k characters
- ▶ or NO_MATCH if there is no such i

"Hamming distance $\leq k$ "

mismatched characters



↪ searching with typos

- ▶ Assume longest common extensions in $T_1 P_2$ can be found in $O(1)$
 - ↪ generalized suffix tree \mathcal{T} has been built
 - ↪ string depths of all internal nodes have been computed
 - ↪ constant-time LCA data structure for \mathcal{T} has been built

Clicker Question



What is the Hamming distance between heart and beard?

2

sli.do/comp526

Click on "Polls" tab

A large, light pink arrow pointing to the right is located at the bottom right of the slide, containing the text "Click on 'Polls' tab".

Clicker Question



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

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

sli.do/comp526

Click on "Polls" tab

Clicker Question



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

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

sli.do/comp526

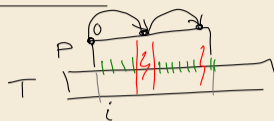
Click on "Polls" tab

Kangaroo Algorithm for approximate matching

easy in $O(n \cdot m)$



```
1 procedure kMismatch( $T[0..n - 1], P[0..m - 1]$ )
2   // build LCE data structure
3   for  $i := 0, \dots, n - m - 1$  do
4     mismatches := 0;  $t := i$ ;  $p := 0$ 
5     while mismatches  $\leq k \wedge p < m$  do
6        $\ell := \text{LCE}(t, p)$  // jump over matching part
7        $t := t + \ell + 1$ ;  $p := p + \ell + 1$ 
8       mismatches := mismatches + 1
9     if  $p == m$  then
10      return  $i$ 
```



► **Analysis:** $\Theta(n + m)$ preprocessing + $O(n \cdot k)$ matching

↪ very efficient for small k

► State of the art

► $O\left(n \frac{k^2 \log k}{m}\right)$ possible with complicated algorithms

► extensions for edit distance $\leq k$ possible

Application 6: Matching with wildcards

- ▶ Allow a wildcard character in pattern

stands for arbitrary (single) character

unit*	<i>P</i>
in_unit5_we_will	<i>T</i>

- ▶ 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*	P
in_unit5_we_will	T

- ▶ 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

👍 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

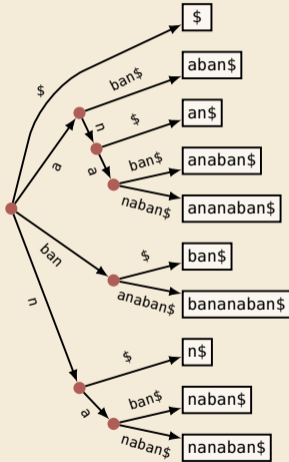


∪ ∑ = ASCII 6 = 128

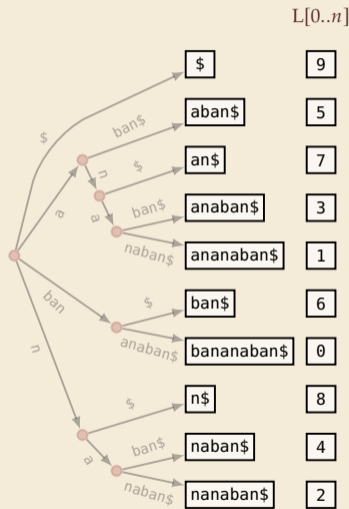
6.5 Suffix Arrays

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

► **Example:** $P = ana$

↪ $L[0..n]$ is called *suffix array*:

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

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

Clicker Question

Check all correct statements about suffix array $L[0..n]$ and suffix tree \mathcal{T} of text $T[0..n]$.



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

sli.do/comp526

Click on "Polls" tab

Clicker Question

Check all correct statements about suffix array $L[0..n]$ and suffix tree \mathcal{T} of text $T[0..n]$.



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

sli.do/comp526

Click on "Polls" tab

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 do better!

Fat-pivot radix quicksort – Example *(corrected version)*

she

sells

seashells

by

the

sea

shore

the

shells

she

sells

are

surely

seashells

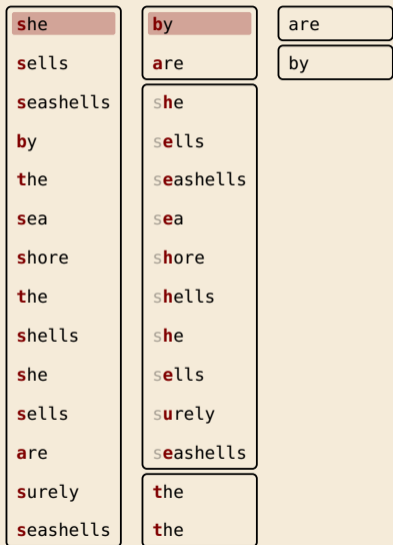
Fat-pivot radix quicksort – Example

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



Fat-pivot radix quicksort – Example



Fat-pivot radix quicksort – Example

she	by	are
sells	are	by
seashells	she	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



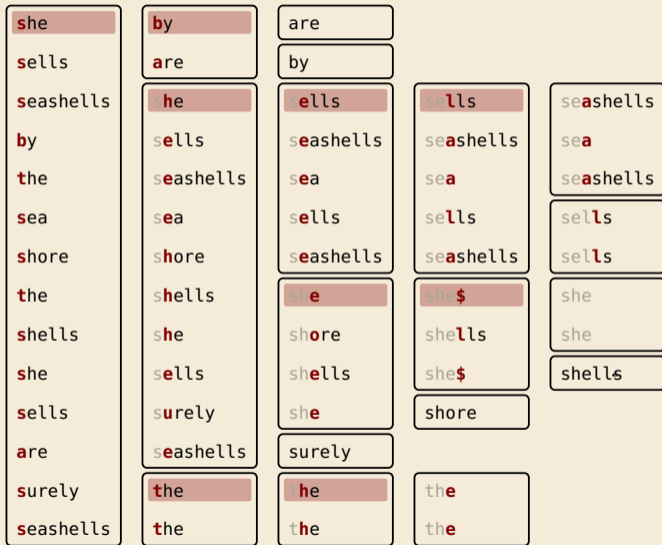
Fat-pivot radix quicksort – Example



Fat-pivot radix quicksort – Example



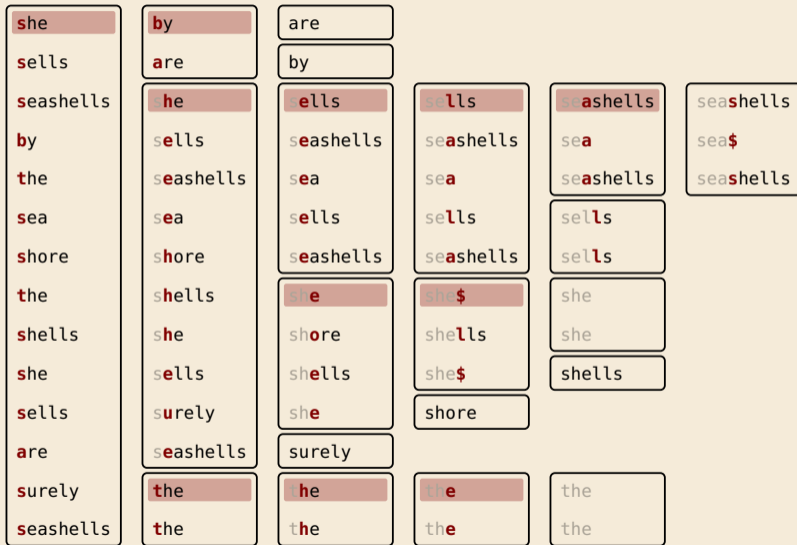
Fat-pivot radix quicksort – Example



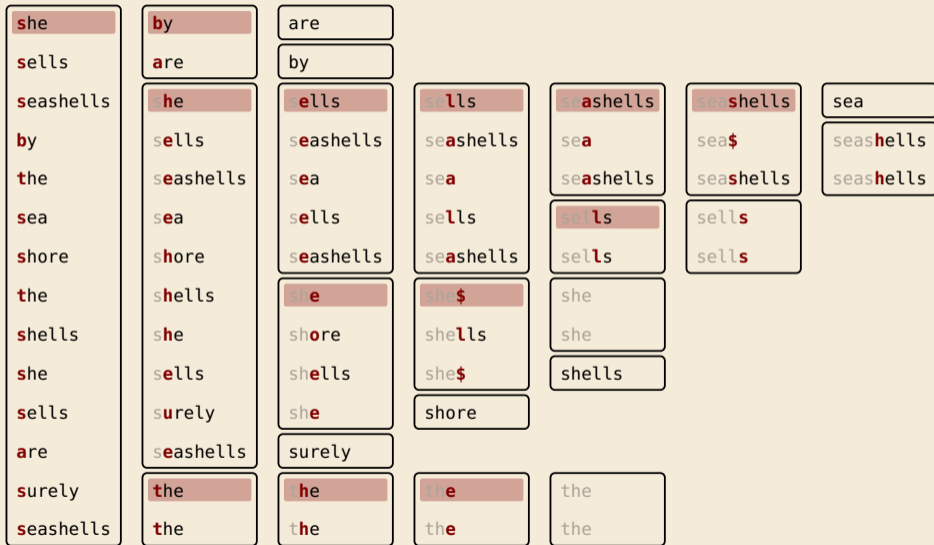
Fat-pivot radix quicksort – Example



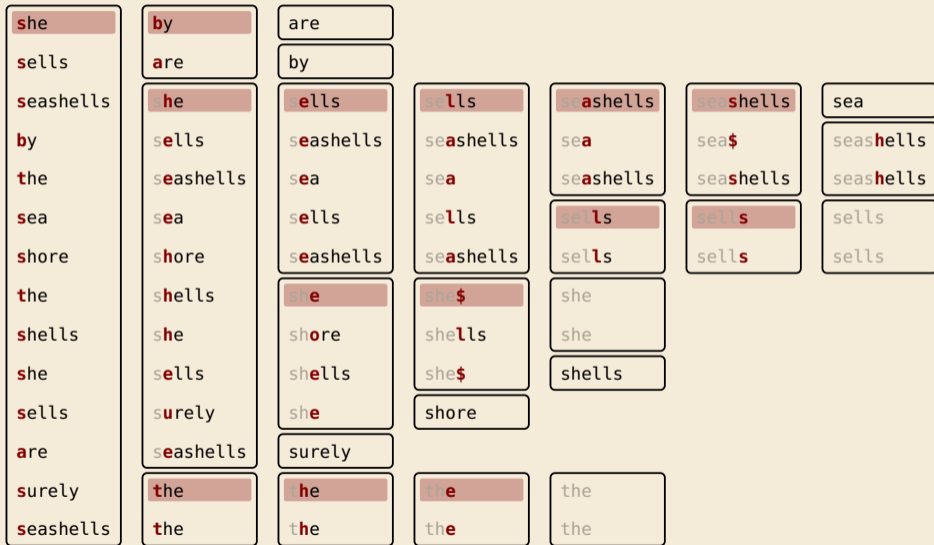
Fat-pivot radix quicksort – Example



Fat-pivot radix quicksort – Example



Fat-pivot radix quicksort – Example



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 $d = 0$)
- ↪ 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

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

for random strings

choice of pivot
and random strings

👍 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

w.c. $T = a a a c a \dots \$$

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

↪ can show: $\sim 2 \ln(2) \cdot n \lg n \approx 1.39n \lg n$ character comparisons on average for random strings

👍 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

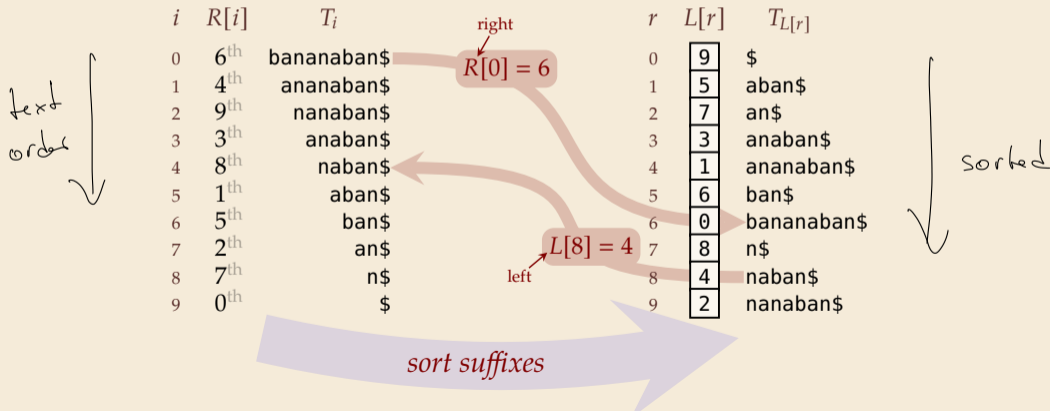
but we can do $O(n)$ time worst case!

6.6 Linear-Time Suffix Sorting

Inverse suffix array: going left & right

► to understand the fastest algorithm, it is helpful to define the *inverse suffix array*:

- $R[i] = r \iff L[r] = i$ $L = \text{leaf array}$
- \iff there are r suffixes that come before T_i in sorted order
- $\iff T_i$ has (0-based) *rank* $r \rightsquigarrow$ call $R[0..n]$ the rank array



Linear-time suffix sorting

DC3 / Skew algorithm

1. Compute rank array $R_{1,2}$ for suffixes T_i starting at $i \not\equiv 0 \pmod{3}$ *recursively*. *not a multiple of 3* $\frac{2}{3} \surd$
2. Induce rank array R_3 for suffixes $T_0, T_3, T_6, T_9, \dots$ from $R_{1,2}$.
3. Merge $R_{1,2}$ and R_0 using $R_{1,2}$.
 \rightsquigarrow rank array R for entire input

Linear-time suffix sorting

DC3 / Skew algorithm

1. Compute rank array $R_{1,2}$ for suffixes T_i starting at $i \not\equiv 0 \pmod{3}$ *not a multiple of 3* recursively.
2. Induce rank array R_3 for suffixes $T_0, T_3, T_6, T_9, \dots$ from $R_{1,2}$.
3. Merge $R_{1,2}$ and R_3 using $R_{1,2}$.
 \rightsquigarrow rank array R for entire input

► We will show that steps 2. and 3. take $\Theta(n)$ time .

\rightsquigarrow Total complexity is $\underbrace{n}_{\text{top level}} + \underbrace{\frac{2}{3}n}_{\text{second level}} + \left(\frac{2}{3}\right)^2 n + \left(\frac{2}{3}\right)^3 n + \dots \leq n \cdot \sum_{i \geq 0} \left(\frac{2}{3}\right)^i = 3n = \underline{\Theta(n)}$

Linear-time suffix sorting

DC3 / Skew algorithm

1. Compute rank array $R_{1,2}$ for suffixes T_i starting at $i \not\equiv 0 \pmod{3}$ *not a multiple of 3* recursively.
2. Induce rank array R_3 for suffixes $T_0, T_3, T_6, T_9, \dots$ from $R_{1,2}$.
3. Merge $R_{1,2}$ and R_0 using $R_{1,2}$.
 \rightsquigarrow rank array R for entire input

► We will show that steps 2. and 3. take $\Theta(n)$ time

\rightsquigarrow Total complexity is $n + \frac{2}{3}n + \left(\frac{2}{3}\right)^2 n + \left(\frac{2}{3}\right)^3 n + \dots \leq n \cdot \sum_{i \geq 0} \left(\frac{2}{3}\right)^i = 3n = \Theta(n)$

► **Note:** L can easily be computed from R in one pass, and vice versa.

\rightsquigarrow Can use whichever is more convenient.

DC3 / Skew algorithm – Step 2: Inducing ranks

- ▶ **Assume:** rank array $R_{1,2}$ known:

- ▶ $R_{1,2}[i] = \begin{cases} \text{rank of } T_i \text{ among } T_1, T_2, T_4, T_5, T_7, T_8, \dots & \text{for } i = 1, 2, 4, 5, 7, 8, \dots \\ \text{undefined} & \text{for } i = 0, 3, 6, 9, \dots \end{cases}$

- ▶ **Task:** sort the suffixes $T_0, T_3, T_6, T_9, \dots$ in linear time (!)

DC3 / Skew algorithm – Step 2: Inducing ranks

- ▶ **Assume:** rank array $R_{1,2}$ known:

$$\text{▶ } R_{1,2}[i] = \begin{cases} \text{rank of } T_i \text{ among } T_1, T_2, T_4, T_5, T_7, T_8, \dots & \text{for } i = 1, 2, 4, 5, 7, 8, \dots \\ \text{undefined} & \text{for } i = 0, 3, 6, 9, \dots \end{cases}$$

- ▶ **Task:** sort the suffixes $T_0, T_3, T_6, T_9, \dots$ in linear time (!)

- ▶ Suppose we want to compare T_0 and T_3 .

$$T_0 = a \overbrace{\quad T_1 \quad}$$

$$T_3 = c \overbrace{\quad T_4 \quad}$$

- ▶ Characterwise comparisons too expensive
- ▶ but: after removing first character, we obtain T_1 and T_4
- ▶ these two can be compared in *constant time* by comparing $\underline{R_{1,2}[1]}$ and $\underline{R_{1,2}[4]}$!

⇒ T_0 comes before T_3 in lexicographic order
iff pair $(T[0], R_{1,2}[1])$ comes before pair $(T[3], R_{1,2}[4])$ in lexicographic order

DC3 / Skew algorithm – Inducing ranks example

$T = \text{hannahbansbananasman}\$\$\$$

(append 3 \$ markers)

T_0 hannahbansbananasman\$\\$\\$
 T_3 nahbansbananasman\$\\$\\$
 T_6 bansbananasman\$\\$\\$
 T_9 sbananasman\$\\$\\$
 T_{12} nanasman\$\\$\\$
 T_{15} asman\$\\$\\$
 T_{18} an\$\\$\\$
 T_{21} \$\$

T_1	annahbansbananasman\$\\$\\$	$R_{1,2}[22] = 0$	T_{22}	\$
T_2	nahbansbananasman\$\\$\\$	$R_{1,2}[20] = 1$	T_{20}	\$\$
T_4	ahbansbananasman\$\\$\\$	$R_{1,2}[4] = 2$	T_4	ahbansbananasman\$\\$\\$
T_5	hbansbananasman\$\\$\\$	$R_{1,2}[11] = 3$	T_{11}	anasman\$\\$\\$
T_7	ansbananasman\$\\$\\$	$R_{1,2}[13] = 4$	T_{13}	anasman\$\\$\\$
T_8	nsbananasman\$\\$\\$	$R_{1,2}[1] = 5$	T_1	annahbansbananasman\$\\$\\$
T_{10}	bananasman\$\\$\\$	$R_{1,2}[7] = 6$	T_7	ansbananasman\$\\$\\$
T_{11}	anasman\$\\$\\$	$R_{1,2}[10] = 7$	T_{10}	bananasman\$\\$\\$
T_{13}	anasman\$\\$\\$	$R_{1,2}[5] = 8$	T_5	hbansbananasman\$\\$\\$
T_{14}	nasman\$\\$\\$	$R_{1,2}[17] = 9$	T_{17}	man\$\\$\\$
T_{16}	sman\$\\$\\$	$R_{1,2}[19] = 10$	T_{19}	n\$\\$\\$
T_{17}	man\$\\$\\$	$R_{1,2}[14] = 11$	T_{14}	nasman\$\\$\\$
T_{19}	n\$\\$\\$	$R_{1,2}[2] = 12$	T_2	nahbansbananasman\$\\$\\$
T_{20}	\$\$	$R_{1,2}[8] = 13$	T_8	nsbananasman\$\\$\\$
T_{22}	\$	$R_{1,2}[16] = 14$	T_{16}	sman\$\\$\\$

$R_{1,2}$ (known)

DC3 / Skew algorithm – Inducing ranks example

$T = \text{hannahbansbananasman}\$\$\$$

(append 3 \$ markers)

T_0 hannahbansbananasman\$\\$\\$
 T_3 nahbansbananasman\$\\$\\$
 T_6 bansbananasman\$\\$\\$
 T_9 sbananasman\$\\$\\$
 T_{12} nanasman\$\\$\\$
 T_{15} asman\$\\$\\$
 T_{18} an\$\\$\\$
 T_{21} \$\$

$\text{sman}\$\$\$ = T_{16}$

T_0 h05
 T_3 n02
 T_6 b06
 T_9 s07
 T_{12} n04
 T_{15} a14
 T_{18} a10
 T_{21} \$00

$R_{1,2}[16] = 14$

T_1	annahbansbananasman\$\\$\\$	$R_{1,2}[22] = 0$	T_{22}	\$
T_2	nnaahbansbananasman\$\\$\\$	$R_{1,2}[20] = 1$	T_{20}	\$\$
T_4	ahbansbananasman\$\\$\\$	$R_{1,2}[4] = 2$	T_4	ahbansbananasman\$\\$\\$
T_5	hbansbananasman\$\\$\\$	$R_{1,2}[11] = 3$	T_{11}	anasman\$\\$\\$
T_7	ansbananasman\$\\$\\$	$R_{1,2}[13] = 4$	T_{13}	anasman\$\\$\\$
T_8	nsbananasman\$\\$\\$	$R_{1,2}[1] = 5$	T_1	annahbansbananasman\$\\$\\$
T_{10}	bananasman\$\\$\\$	$R_{1,2}[7] = 6$	T_7	ansbananasman\$\\$\\$
T_{11}	anasman\$\\$\\$	$R_{1,2}[10] = 7$	T_{10}	bananasman\$\\$\\$
T_{13}	anasman\$\\$\\$	$R_{1,2}[5] = 8$	T_5	hbansbananasman\$\\$\\$
T_{14}	nasman\$\\$\\$	$R_{1,2}[17] = 9$	T_{17}	man\$\\$\\$
T_{16}	sman\$\\$\\$	$R_{1,2}[19] = 10$	T_{19}	n\$\\$\\$
T_{17}	man\$\\$\\$	$R_{1,2}[14] = 11$	T_{14}	nasman\$\\$\\$
T_{19}	n\$\\$\\$	$R_{1,2}[2] = 12$	T_2	nnaahbansbananasman\$\\$\\$
T_{20}	\$\$	$R_{1,2}[8] = 13$	T_8	nsbananasman\$\\$\\$
T_{22}	\$	$R_{1,2}[16] = 14$	T_{16}	sman\$\\$\\$

$R_{1,2}$ (known)

DC3 / Skew algorithm – Inducing ranks example

$T = \text{hannahbansbananasman}\$ \$ \$$

(append 3 \$ markers)

T_0 hannahbansbananasman\$\$\$
 T_3 nahbansbananasman\$\$\$
 T_6 bansbananasman\$\$\$
 T_9 sbananasman\$\$\$
 T_{12} nanasman\$\$\$
 T_{15} asman\$\$\$
 T_{18} an\$\$\$
 T_{21} \$\$

smans\$\$\$ = T_{16}

T_0 h05
 T_3 n02
 T_6 b06
 T_9 s07
 T_{12} n04
 T_{15} a14
 T_{18} a10
 T_{21} \$00

$R_{1,2}[16] = 14$

T_1	annahbansbananasman\$\$\$	$R_{1,2}[22] = 0$	T_{22}	\$
T_2	nahbansbananasman\$\$\$	$R_{1,2}[20] = 1$	T_{20}	\$\$\$
T_4	ahbansbananasman\$\$\$	$R_{1,2}[4] = 2$	T_4	ahbansbananasman\$\$\$
T_5	hbansbananasman\$\$\$	$R_{1,2}[11] = 3$	T_{11}	anasman\$\$\$
T_7	ansbananasman\$\$\$	$R_{1,2}[13] = 4$	T_{13}	anasman\$\$\$
T_8	nsbananasman\$\$\$	$R_{1,2}[1] = 5$	T_1	annahbansbananasman\$\$\$
T_{10}	bananasman\$\$\$	$R_{1,2}[7] = 6$	T_7	ansbananasman\$\$\$
T_{11}	anasman\$\$\$	$R_{1,2}[10] = 7$	T_{10}	bananasman\$\$\$
T_{13}	anasman\$\$\$	$R_{1,2}[5] = 8$	T_5	hbansbananasman\$\$\$
T_{14}	nasman\$\$\$	$R_{1,2}[17] = 9$	T_{17}	man\$\$\$
T_{16}	smans\$\$\$	$R_{1,2}[19] = 10$	T_{19}	n\$\$\$
T_{17}	mans\$\$\$	$R_{1,2}[14] = 11$	T_{14}	nasman\$\$\$
T_{19}	n\$\$\$	$R_{1,2}[2] = 12$	T_2	nahbansbananasman\$\$\$
T_{20}	\$\$\$	$R_{1,2}[8] = 13$	T_8	nsbananasman\$\$\$
T_{22}	\$	$R_{1,2}[16] = 14$	T_{16}	smans\$\$\$

$R_{1,2}$ (known)



T_{21} \$00 \rightsquigarrow $R_0[21] = 0$
 T_{18} a10 \rightsquigarrow $R_0[18] = 1$
 T_{15} a14 \rightsquigarrow $R_0[15] = 2$
 T_6 b06 \rightsquigarrow $R_0[6] = 3$
 T_0 h05 \rightsquigarrow $R_0[0] = 4$
 T_3 n02 \rightsquigarrow $R_0[3] = 5$
 T_{12} n04 \rightsquigarrow $R_0[12] = 6$
 T_9 s07 \rightsquigarrow $R_0[9] = 7$

DC3 / Skew algorithm – Inducing ranks example

$T = \text{hannahbansbananasman}\$\$\$$

(append 3 \$ markers)

T_0 hannahbansbananasman\$\$\$
 T_3 nahbansbananasman\$\$\$
 T_6 bansbananasman\$\$\$
 T_9 sbananasman\$\$\$
 T_{12} nanasman\$\$\$
 T_{15} asman\$\$\$
 T_{18} an\$\$\$
 T_{21} \$\$

smans\$\$\$ = T_{16}

T_0 h05
 T_3 n02
 T_6 b06
 T_9 s07
 T_{12} n04
 T_{15} a14
 T_{18} a10
 T_{21} \$00

$R_{1,2}[16] = 14$

T_1	annahbansbananasman\$\$\$	$R_{1,2}[22] = 0$	T_{22}	\$
T_2	nahbansbananasman\$\$\$	$R_{1,2}[20] = 1$	T_{20}	\$\$\$
T_4	ahbansbananasman\$\$\$	$R_{1,2}[4] = 2$	T_4	ahbansbananasman\$\$\$
T_5	hbansbananasman\$\$\$	$R_{1,2}[11] = 3$	T_{11}	anasman\$\$\$
T_7	ansbananasman\$\$\$	$R_{1,2}[13] = 4$	T_{13}	anasman\$\$\$
T_8	nsbananasman\$\$\$	$R_{1,2}[1] = 5$	T_1	annahbansbananasman\$\$\$
T_{10}	bananasman\$\$\$	$R_{1,2}[7] = 6$	T_7	ansbananasman\$\$\$
T_{11}	anasman\$\$\$	$R_{1,2}[10] = 7$	T_{10}	bananasman\$\$\$
T_{13}	anasman\$\$\$	$R_{1,2}[5] = 8$	T_5	hbansbananasman\$\$\$
T_{14}	nasman\$\$\$	$R_{1,2}[17] = 9$	T_{17}	man\$\$\$
T_{16}	smans\$\$\$	$R_{1,2}[19] = 10$	T_{19}	n\$\$\$
T_{17}	mans\$\$\$	$R_{1,2}[14] = 11$	T_{14}	nasman\$\$\$
T_{19}	n\$\$\$	$R_{1,2}[2] = 12$	T_2	nahbansbananasman\$\$\$
T_{20}	\$\$\$	$R_{1,2}[8] = 13$	T_8	nsbananasman\$\$\$
T_{22}	\$	$R_{1,2}[16] = 14$	T_{16}	smans\$\$\$

$R_{1,2}$ (known)



T_{21}	\$00	\rightsquigarrow	$R_0[21] = 0$
T_{18}	a10	\rightsquigarrow	$R_0[18] = 1$
T_{15}	a14	\rightsquigarrow	$R_0[15] = 2$
T_6	b06	\rightsquigarrow	$R_0[6] = 3$
T_0	h05	\rightsquigarrow	$R_0[0] = 4$
T_3	n02	\rightsquigarrow	$R_0[3] = 5$
T_{12}	n04	\rightsquigarrow	$R_0[12] = 6$
T_9	s07	\rightsquigarrow	$R_0[9] = 7$

R_0

DC3 / Skew algorithm – Inducing ranks example

$T = \text{hannahbansbananasman}\$\$\$$

(append 3 \$ markers)

T_0 hannahbansbananasman\$\$\$
 T_3 nahbansbananasman\$\$\$
 T_6 bansbananasman\$\$\$
 T_9 sbananasman\$\$\$
 T_{12} nanasman\$\$\$
 T_{15} asman\$\$\$
 T_{18} an\$\$\$
 T_{21} \$\$

smans\$\$\$ = T_{16}

T_0 h05
 T_3 n02
 T_6 b06
 T_9 s07
 T_{12} n04
 T_{15} a14
 T_{18} a10
 T_{21} \$00

$R_{1,2}[16] = 14$

T_1	annahbansbananasman\$\$\$	$R_{1,2}[22] = 0$	T_{22}	\$
T_2	nahbansbananasman\$\$\$	$R_{1,2}[20] = 1$	T_{20}	\$\$\$
T_4	ahbansbananasman\$\$\$	$R_{1,2}[4] = 2$	T_4	ahbansbananasman\$\$\$
T_5	hbansbananasman\$\$\$	$R_{1,2}[11] = 3$	T_{11}	anasman\$\$\$
T_7	ansbananasman\$\$\$	$R_{1,2}[13] = 4$	T_{13}	anasman\$\$\$
T_8	nsbananasman\$\$\$	$R_{1,2}[1] = 5$	T_1	annahbansbananasman\$\$\$
T_{10}	bananasman\$\$\$	$R_{1,2}[7] = 6$	T_7	ansbananasman\$\$\$
T_{11}	anasman\$\$\$	$R_{1,2}[10] = 7$	T_{10}	bananasman\$\$\$
T_{13}	anasman\$\$\$	$R_{1,2}[5] = 8$	T_5	hbansbananasman\$\$\$
T_{14}	nasman\$\$\$	$R_{1,2}[17] = 9$	T_{17}	man\$\$\$
T_{16}	smans\$\$\$	$R_{1,2}[19] = 10$	T_{19}	n\$\$\$
T_{17}	mans\$\$\$	$R_{1,2}[14] = 11$	T_{14}	nasman\$\$\$
T_{19}	n\$\$\$	$R_{1,2}[2] = 12$	T_2	nahbansbananasman\$\$\$
T_{20}	\$\$\$	$R_{1,2}[8] = 13$	T_8	nsbananasman\$\$\$
T_{22}	\$	$R_{1,2}[16] = 14$	T_{16}	smans\$\$\$

$R_{1,2}$ (known)



T_{21}	\$00	\rightsquigarrow	$R_0[21] = 0$
T_{18}	a10	\rightsquigarrow	$R_0[18] = 1$
T_{15}	a14	\rightsquigarrow	$R_0[15] = 2$
T_6	b06	\rightsquigarrow	$R_0[6] = 3$
T_0	h05	\rightsquigarrow	$R_0[0] = 4$
T_3	n02	\rightsquigarrow	$R_0[3] = 5$
T_{12}	n04	\rightsquigarrow	$R_0[12] = 6$
T_9	s07	\rightsquigarrow	$R_0[9] = 7$

R_0

► sorting of pairs doable in $O(n)$ time by 2 iterations of counting sort

\rightsquigarrow Obtain R_0 in $O(n)$ time

Clicker Question

Recap: Check all correct statements about suffix array $L[0..n]$, inverse suffix array $R[0..n]$, and suffix tree \mathcal{T} of text T .



- A** L lists the leaves of \mathcal{T} in left-to-right order.
- B** R lists the leaves of \mathcal{T} in right-to-left order.
- C** R lists starting indices of suffixes in lexicographic order.
- D** L lists starting indices of suffixes in lexicographic order.
- E** $L[r] = i$ iff $R[i] = r$
- F** L stands for leaf
- G** L stands for left
- H** R stands for rank
- I** R stands for right

sli.do/comp526

Click on "Polls" tab

Clicker Question

Recap: Check all correct statements about suffix array $L[0..n]$, inverse suffix array $R[0..n]$, and suffix tree \mathcal{T} of text T .



- A** L lists the leaves of \mathcal{T} in left-to-right order. ✓
- B** ~~R lists the leaves of \mathcal{T} in right to left order.~~
- C** ~~R lists starting indices of suffixes in lexicographic order.~~
- D** L lists starting indices of suffixes in lexicographic order. ✓
- E** $L[r] = i$ iff $R[i] = r$ ✓
- F** L stands for leaf ✓
- G** L stands for left ✓
- H** R stands for rank ✓
- I** R stands for right ✓

sli.do/comp526

Click on "Polls" tab

DC3 / Skew algorithm – Step 3: Merging

T_{21} \$\$
 T_{18} an\$\$\$
 T_{15} asman\$\$\$
 T_6 bansbananasman\$\$\$
 T_0 hannahbansbananasman\$\$\$
 T_3 nahbansbananasman\$\$\$
 T_{12} nanasman\$\$\$
 T_9 sbananasman\$\$\$

T_{22} \$
 T_{20} \$\$\$
 T_4 ahbansbananasman\$\$\$
 T_{11} ananasman\$\$\$
 T_{13} anasman\$\$\$
 T_1 annahbansbananasman\$\$\$
 T_7 ansbananasman\$\$\$
 T_{10} bananasman\$\$\$
 T_5 hbansbananasman\$\$\$
 T_{17} man\$\$\$
 T_{19} n\$\$\$
 T_{14} nasman\$\$\$
 T_2 nnahbansbananasman\$\$\$
 T_8 nsbananasman\$\$\$
 T_{16} sman\$\$\$

► Have:

- sorted 1,2-list:

$T_1, T_2, T_4, T_5, T_7, T_8, T_{10}, T_{11}, \dots$

- sorted 0-list:

$T_0, T_3, T_6, T_9, \dots$

► Task: Merge them!

- use standard merging method from Mergesort
- but speed up comparisons using $R_{1,2}$

DC3 / Skew algorithm – Step 3: Merging

T_{21} \$\$
 T_{18} an\$\$\$
 T_{15} asman\$\$\$
 T_6 bansbananasman\$\$\$
 T_0 hannahbansbananasman\$\$\$
 T_3 nahbansbananasman\$\$\$
 T_{12} nanasman\$\$\$
 T_9 sbanasman\$\$\$

T_{22} \$
 T_{20} \$\$\$
 T_4 ahbansbananasman\$\$\$
 T_{11} anasman\$\$\$
 T_{13} anasman\$\$\$
 T_1 annahbansbananasman\$\$\$
 T_7 ansbananasman\$\$\$
 T_{10} bananasman\$\$\$
 T_5 hbansbananasman\$\$\$
 T_{17} man\$\$\$
 T_{19} n\$\$\$
 T_{14} nasman\$\$\$
 T_2 nnahbansbananasman\$\$\$
 T_8 nsbananasman\$\$\$
 T_{16} sman\$\$\$

T_{22} \$
 T_{21} \$\$
 T_{20} \$\$\$
 T_4 ahbansbananasman\$\$\$
 T_{18} an\$\$\$

► Have:

- sorted 1,2-list:

$T_1, T_2, T_4, T_5, T_7, T_8, T_{10}, T_{11}, \dots$

- sorted 0-list:

$T_0, T_3, T_6, T_9, \dots$

► Task: Merge them!

- use standard merging method from Mergesort
- but speed up comparisons using $R_{1,2}$

DC3 / Skew algorithm – Step 3: Merging

T_{21} \$\$
 T_{18} an\$\$\$
 T_{15} asman\$\$\$
 T_6 bansbananasman\$\$\$
 T_0 hannahbansbananasman\$\$\$
 T_3 nahbansbananasman\$\$\$
 T_{12} nanasman\$\$\$
 T_9 sbanasman\$\$\$

T_{22} \$
 T_{20} \$\$\$
 T_4 ahbansbananasman\$\$\$
 T_{11} ananasman\$\$\$
 T_{13} anasman\$\$\$
 T_1 annahbansbananasman\$\$\$
 T_7 ansbananasman\$\$\$
 T_{10} bananasman\$\$\$
 T_5 hbansbananasman\$\$\$
 T_{17} man\$\$\$
 T_{19} n\$\$\$
 T_{14} nasman\$\$\$
 T_2 nnahbansbananasman\$\$\$
 T_8 nsbananasman\$\$\$
 T_{16} sman\$\$\$

T_{22} \$
 T_{21} \$\$
 T_{20} \$\$\$
 T_4 ahbansbananasman\$\$\$
 T_{18} an\$\$\$

► Have:

- sorted 1,2-list:

$T_1, T_2, T_4, T_5, T_7, T_8, T_{10}, T_{11}, \dots$

- sorted 0-list:

$T_0, T_3, T_6, T_9, \dots$

► Task: Merge them!

- use standard merging method from Mergesort
- but speed up comparisons using $R_{1,2}$

Compare T_{15} to T_{11}

Idea: try same trick as before

$T_{15} = \text{asman}$$$
= asman$$$
= a T_{16}
 $T_{11} = \text{ananasman}$$$
= ananasman$$$
= a $T_{12}$$$

DC3 / Skew algorithm – Step 3: Merging

T_{21} \$\$
 T_{18} an\$\$\$
 T_{15} asman\$\$\$
 T_6 bansbananasman\$\$\$
 T_0 hannahbansbananasman\$\$\$
 T_3 nahbansbananasman\$\$\$
 T_{12} nanasman\$\$\$
 T_9 sbananasman\$\$\$

T_{22} \$
 T_{20} \$\$\$
 T_4 ahbansbananasman\$\$\$
 T_{11} ananasman\$\$\$
 T_{13} anasman\$\$\$
 T_1 annahbansbananasman\$\$\$
 T_7 ansbananasman\$\$\$
 T_{10} bananasman\$\$\$
 T_5 hbansbananasman\$\$\$
 T_{17} man\$\$\$
 T_{19} n\$\$\$
 T_{14} nasman\$\$\$
 T_2 nnahbansbananasman\$\$\$
 T_8 nsbananasman\$\$\$
 T_{16} sman\$\$\$

T_{22} \$
 T_{21} \$\$
 T_{20} \$\$\$
 T_4 ahbansbananasman\$\$\$
 T_{18} an\$\$\$

► Have:

- sorted 1,2-list:

$T_1, T_2, T_4, T_5, T_7, T_8, T_{10}, T_{11}, \dots$

- sorted 0-list:

$T_0, T_3, T_6, T_9, \dots$

► Task: Merge them!

- use standard merging method from Mergesort
- but speed up comparisons using $R_{1,2}$

Compare T_{15} to T_{11}

Idea: try same trick as before

$T_{15} = \text{asman}$$$$

$= \text{asman}$$$$

$= aT_{16}$

can't compare T_{16}
and T_{12} either!

$T_{11} = \text{ananasman}$$$$

$= \text{ananasman}$$$$

$= aT_{12}$

DC3 / Skew algorithm – Step 3: Merging

T_{21} \$\$
 T_{18} an\$\$\$
 T_{15} asman\$\$\$
 T_6 bansbananasman\$\$\$
 T_0 hannahbansbananasman\$\$\$
 T_3 nahbansbananasman\$\$\$
 T_{12} nanasman\$\$\$
 T_9 sbananasman\$\$\$

T_{22} \$
 T_{20} \$\$\$
 T_4 ahbansbananasman\$\$\$
 T_{11} ananasman\$\$\$
 T_{13} anasman\$\$\$
 T_1 annahbansbananasman\$\$\$
 T_7 ansbananasman\$\$\$
 T_{10} bananasman\$\$\$
 T_5 hbansbananasman\$\$\$
 T_{17} man\$\$\$
 T_{19} n\$\$\$
 T_{14} nasman\$\$\$
 T_2 nnahbansbananasman\$\$\$
 T_8 nsbananasman\$\$\$
 T_{16} sman\$\$\$

T_{22} \$
 T_{21} \$\$
 T_{20} \$\$\$
 T_4 ahbansbananasman\$\$\$
 T_{18} an\$\$\$

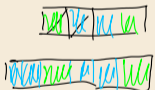
► Have:

- sorted 1,2-list:

$T_1, T_2, T_4, T_5, T_7, T_8, T_{10}, T_{11}, \dots$

- sorted 0-list:

$T_0, T_3, T_6, T_9, \dots$



► Task: Merge them!

- use standard merging method from Mergesort
- but speed up comparisons using $R_{1,2}$

Compare T_{15} to T_{11}

Idea: try same trick as before

$T_{15} =$ asman\$\$\$
 $=$ asman\$\$\$
 $=$ a T_{16}

can't compare T_{16}
and T_{12} either!

$T_{11} =$ ananasman\$\$\$
 $=$ ananasman\$\$\$
 $=$ a T_{12}

↪ Compare T_{16} to T_{12}

$T_{16} =$ sman\$\$\$
 $=$ sman\$\$\$
 $=$ s T_{17}

$T_{12} =$ nanasman\$\$\$
 $=$ aanasman\$\$\$
 $=$ a T_{13}

DC3 / Skew algorithm – Step 3: Merging

T_{21} \$\$
 T_{18} an\$\$\$
 T_{15} asman\$\$\$
 T_6 bansbananasman\$\$\$
 T_0 hannahbansbananasman\$\$\$
 T_3 nahbansbananasman\$\$\$
 T_{12} nanasman\$\$\$
 T_9 sbananasman\$\$\$

T_{22} \$
 T_{20} \$\$\$
 T_4 ahbansbananasman\$\$\$
 T_{11} ananasman\$\$\$
 T_{13} anasman\$\$\$
 T_1 annahbansbananasman\$\$\$
 T_7 ansbananasman\$\$\$
 T_{10} bananasman\$\$\$
 T_5 hbansbananasman\$\$\$
 T_{17} man\$\$\$
 T_{19} n\$\$\$
 T_{14} nasman\$\$\$
 T_2 nnahbansbananasman\$\$\$
 T_8 nsbananasman\$\$\$
 T_{16} sman\$\$\$

T_{22} \$
 T_{21} \$\$
 T_{20} \$\$\$
 T_4 ahbansbananasman\$\$\$
 T_{18} an\$\$\$

► Have:

- sorted 1,2-list:

$T_1, T_2, T_4, T_5, T_7, T_8, T_{10}, T_{11}, \dots$

- sorted 0-list:

$T_0, T_3, T_6, T_9, \dots$

► Task: Merge them!

- use standard merging method from Mergesort
- but speed up comparisons using $R_{1,2}$

Compare T_{15} to T_{11}

Idea: try same trick as before

$T_{15} = \text{asman}$$$
 $= \text{asman}$$$ \quad \text{can't compare } T_{16}$
 $= aT_{16} \quad \text{and } T_{12} \text{ either!}$
 $T_{11} = \text{ananasman}$$$
 $= \text{ananasman}$$$
 $= aT_{12}$$$$

↪ Compare T_{16} to T_{12}

$T_{16} = \text{sman}$$$
 $= \text{sman}$$$ \quad \text{always at most 2 steps}$
 $= sT_{17} \quad \text{then can use } R_{1,2}!$
 $T_{12} = \text{nanasman}$$$
 $= \text{aanasmansman}$$$
 $= aT_{13}$$$$

DC3 / Skew algorithm – Step 3: Merging

T_{21} \$\$
 T_{18} an\$\$\$
 T_{15} asman\$\$\$
 T_6 bansbananasman\$\$\$
 T_0 hannahbansbananasman\$\$\$
 T_3 nahbansbananasman\$\$\$
 T_{12} nanasman\$\$\$
 T_9 sbanasman\$\$\$

T_{22} \$
 T_{20} \$\$\$
 T_4 ahbansbananasman\$\$\$
 T_{11} ananasman\$\$\$
 T_{13} anasman\$\$\$
 T_1 annahbansbananasman\$\$\$
 T_7 ansbananasman\$\$\$
 T_{10} bananasman\$\$\$
 T_5 hbansbananasman\$\$\$
 T_{17} man\$\$\$
 T_{19} n\$\$\$
 T_{14} nasman\$\$\$
 T_2 nnahbansbananasman\$\$\$
 T_8 nsbananasman\$\$\$
 T_{16} sman\$\$\$

T_{22} \$
 T_{21} \$\$
 T_{20} \$\$\$
 T_4 ahbansbananasman\$\$\$
 T_{18} an\$\$\$
 T_c .

► Have:

- sorted 1,2-list:

$T_1, T_2, T_4, T_5, T_7, T_8, T_{10}, T_{11}, \dots$

- sorted 0-list:

$T_0, T_3, T_6, T_9, \dots$

► Task: Merge them!

- use standard merging method from Mergesort
- but speed up comparisons using $R_{1,2}$

↪ $O(n)$ time for merge

Compare T_{15} to T_{11}

Idea: try same trick as before

$T_{15} = \text{asman}$$$
 $= \text{asman}$$$ \quad \text{can't compare } T_{16}$
 $= aT_{16} \quad \text{and } T_{12} \text{ either!}$
 $T_{11} = \text{ananasman}$$$
 $= \text{ananasman}$$$
 $= aT_{12}$$$$

↪ Compare T_{16} to T_{12}

$T_{16} = \text{sman}$$$
 $= \text{sman}$$$ \quad \text{always at most 2 steps}$
 $= sT_{17} \quad \text{then can use } R_{1,2}!$
 $T_{12} = \text{nanasman}$$$
 $= \text{aananasman}$$$
 $= aT_{13}$$$$

DC3 / Skew algorithm – Fix recursive call

- ▶ both step 2. and 3. doable in $O(n)$ time!

DC3 / Skew algorithm – Fix recursive call

- ▶ both step 2. and 3. doable in $O(n)$ time!
- ▶ But: we cheated in 1. step! “compute rank array $R_{1,2}$ recursively”
 - ▶ Taking a *subset* of suffixes is *not* an instance of the same problem!



DC3 / Skew algorithm – Fix recursive call

- ▶ both step 2. and 3. doable in $O(n)$ time!
 - ▶ But: we cheated in 1. step! “compute rank array $R_{1,2}$ recursively”
 - ▶ Taking a *subset* of suffixes is *not* an instance of the same problem!
- ↪ Need a single *string* T' to recurse on, from which we can deduce $R_{1,2}$.



How can we make T' “skip” some suffixes?

DC3 / Skew algorithm – Fix recursive call

▶ both step 2. and 3. doable in $O(n)$ time!

▶ But: we cheated in 1. step! “compute rank array $R_{1,2}$ recursively”

▶ Taking a *subset* of suffixes is *not* an instance of the same problem!

↪ Need a single *string* T' to recurse on, from which we can deduce $R_{1,2}$.



How can we make T' “skip” some suffixes?



redefine alphabet to be *triples of characters* \boxed{abc}

$T = \text{bananaban}\$\$\$\$$
 ↪ $T^\square = \boxed{\text{ban}}\boxed{\text{ana}}\boxed{\text{ban}}\boxed{\$\$\$}$
 $\boxed{\text{ana}}\boxed{\text{ban}}\boxed{\$\$\$}$
 $\boxed{\text{ban}}\boxed{\$\$\$}$
 $\boxed{\$\$\$}$

↪ suffixes of $T^\square \iff T_0, T_3, T_6, T_9, \dots$

$T_1, T_4, T_7, \dots \quad T_2, T_5, T_8, \dots$

▶ $T' = \boxed{T[1..n]^\square \ \$\$\$ T[2..n]^\square \ \$\$\$} \iff T_i \text{ with } i \not\equiv 0 \pmod{3}.$

↪ Can call suffix sorting recursively on T' and map result to $R_{1,2}$

DC3 / Skew algorithm – Fix alphabet explosion

- ▶ Still does not quite work!

DC3 / Skew algorithm – Fix alphabet explosion

- ▶ Still does not quite work!

- ▶ Each recursive step *cubes* σ by using triples!

- ↪ (Eventually) cannot use linear-time sorting anymore!

σ

σ^3

$$(\sigma^3)^3 = \sigma^9$$

2^{30}

DC3 / Skew algorithm – Fix alphabet explosion

- ▶ Still does not quite work!
 - ▶ Each recursive step *cubes* σ by using triples!
 - ↪ (Eventually) cannot use linear-time sorting anymore!

▶ But: Have at most $\frac{2}{3}n$ different triples \boxed{abc} in T' !

↪ Before recursion:

1. Sort all occurring triples. (using counting sort in $O(n)$)
2. Replace them by their *rank* (in Σ).

↪ Maintains $\sigma \leq n$ without affecting order of suffixes.

DC3 / Skew algorithm – Step 3. Example

$$T' = T[1..n) \square \square \square T[2..n) \square \square \square$$

► $T = \text{hannahbansbananasman\$}$

DC3 / Skew algorithm – Step 3. Example

$$T' = T[1..n) \square \square \square T[2..n) \square \square \square$$

► $T = \text{hannahbansbananasman} \square \square \square T_2 = \text{nnaahbansbananasman} \square \square \square$

$T' = \text{annahbansbananasman} \square \square \square \text{nnaahbansbananasman} \square \square \square$

DC3 / Skew algorithm – Step 3. Example

$$T' = T[1..n) \square \square \square T[2..n) \square \square \square$$

► $T = \text{hannahbansbananasman\$}$ $T_2 = \text{nna hba nsb ana nas man\$}$

$T' = \text{ann ahb ans ban ana sma n\$\$} \square \square \square \text{nna hba nsb ana nas man \$\$\$}$

► Occurring triples:

$\text{ann ahb ans ban ana sma n\$\$} \square \square \square \text{nna hba nsb} \quad \text{nas man}$

DC3 / Skew algorithm – Step 3. Example

$$T' = T[1..n) \square \square \square T[2..n) \square \square \square$$

- ▶ $T = \text{hannahbansbananasman\$}$ $T_2 = \text{nahbansbananasman\$}$
 $T' = \text{annahbansbananasman\$ \$\$ \$\$ \$\$ nnahbansbananasman\$ \$\$ \$\$}$

- ▶ Occurring triples:

annahbansbananasman\\$ \\$\\$ \\$\\$ nnahbansbananasman\\$ \\$\\$ \\$\\$ nasman

- ▶ Sorted triples with ranks:

Rank	00	01	02	03	04	05	06	07	08	09	10	11	12
Triple	\$\$\$	ahb	ana	ann	ans	ban	hba	man	n\$\$	nas	nna	nsb	sma

DC3 / Skew algorithm – Step 3. Example

$$T' = T[1..n) \square \square \square T[2..n) \square \square \square$$

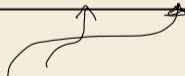
▶ $T = \text{hannahbansbananasman\$}$ $T_2 = \text{nnaahbansbananasman\$}$
 $T' = \text{annahbansbananasman\$ \$\$ \$\$ nnaahbansbananasman \$\$ \$\$}$

▶ Occurring triples:

$\text{annahbansbananasman\$ \$\$ \$\$ nnaahbansbananasman \$\$ \$\$}$

▶ Sorted triples with ranks:

Rank	00	01	02	03	04	05	06	07	08	09	10	11	12
Triple	\$\$\$	ahb	ana	ann	ans	ban	hba	man	n\$\$	nas	nna	nsb	sma



▶ $T' = \text{annahbansbananasman\$ \$\$ \$\$ nnaahbansbananasman \$\$ \$\$}$
 $T'' = \text{03 01 04 05 02 12 08 00 10 06 11 02 09 07 00}$

Suffix array – Discussion

👍 sleek data structure compared to suffix tree

👍 simple and fast $O(n \log n)$ construction (on random strings)

👍 more involved but optimal $O(n)$ construction

👍 supports efficient string matching

👎 string matching takes $O(m \log n)$, not optimal $O(m)$

👎 Cannot use more advanced suffix tree features
e. g., for longest repeated substrings



6.7 The LCP Array

Clicker Question



Which feature of suffix **trees** did we use to find the *length* of a longest repeated substring?

- A** order of leaves
- B** path label of internal nodes
- C** string depth of internal nodes
- D** constant-time traversal to child nodes
- E** constant-time traversal to parent nodes
- F** constant-time traversal to leftmost leaf in subtree

sli.do/comp526

Click on "Polls" tab

Clicker Question



Which feature of suffix **trees** did we use to find the *length* of a longest repeated substring?

- ~~A order of leaves~~
- ~~B path label of internal nodes~~
- C string depth of internal nodes ✓
- ~~D constant time traversal to child nodes~~
- ~~E constant time traversal to parent nodes~~
- ~~F constant time traversal to leftmost leaf in subtree~~

sli.do/comp526

Click on "Polls" tab

String depths of internal nodes

- ▶ Recall algorithm for longest repeated substring in **suffix tree**

1. Compute *string depth* of nodes
2. Find *path label* to node with maximal string depth

- ▶ Can we do this using **suffix arrays**?

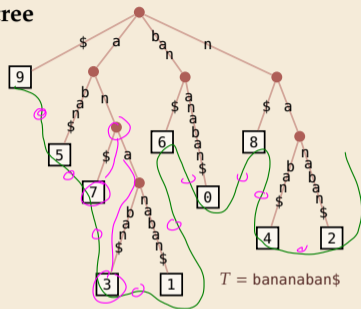
- ▶ Yes, by **enhancing** the suffix array with the **LCP array!**

$LCP[1..n]$

$$LCP[r] = LCP(T_{L[r]}, T_{L[r-1]})$$

length of longest common prefix of suffixes of rank r and $r - 1$

↪ longest repeated substring = find maximum in $LCP[1..n]$



LCP array and internal nodes

$L[0..n]$

9

5

7

3

1

6

0

8

4

2

LCP array and internal nodes

	L[0..n]
\$	9
aban\$	5
an\$	7
anaban\$	3
ananaban\$	1
ban\$	6
bananaban\$	0
n\$	8
naban\$	4
nanaban\$	2

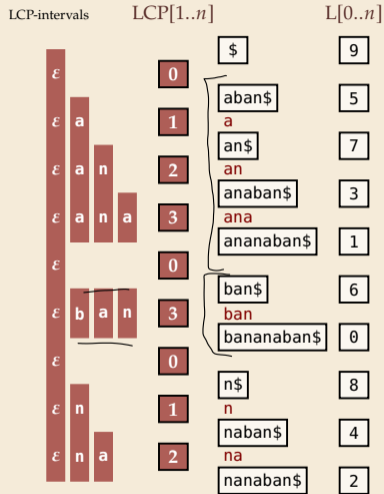
LCP array and internal nodes

	L[0..n]
\$	9
aban\$	5
a	
an\$	7
an	
anaban\$	3
ana	
ananaban\$	1
ban\$	6
ban	
bananaban\$	0
n\$	8
n	
naban\$	4
na	
nanaban\$	2

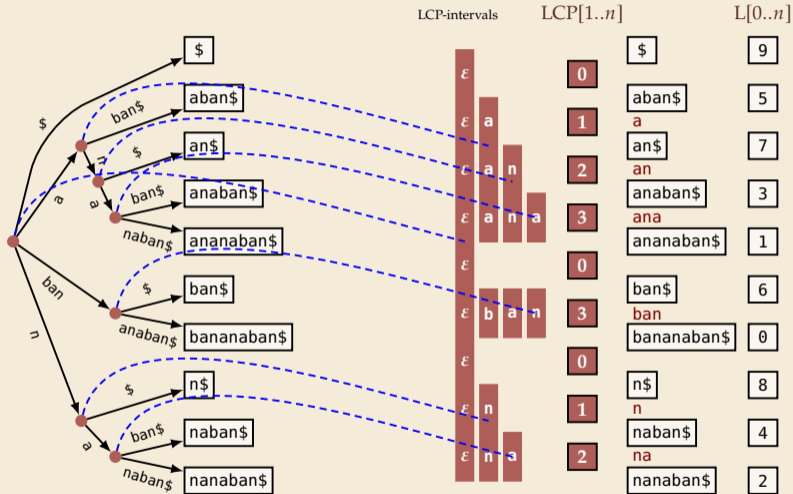
LCP array and internal nodes

	LCP[1..n]	L[0..n]
	\$	9
0	aban\$	5
1	a	7
2	an\$	3
3	an	1
0	anaban\$	6
3	ana	0
0	ananaban\$	8
3	ban\$	4
0	ban	2
0	banaban\$	
1	n\$	
1	n	
2	naban\$	
2	na	
	nanaban\$	

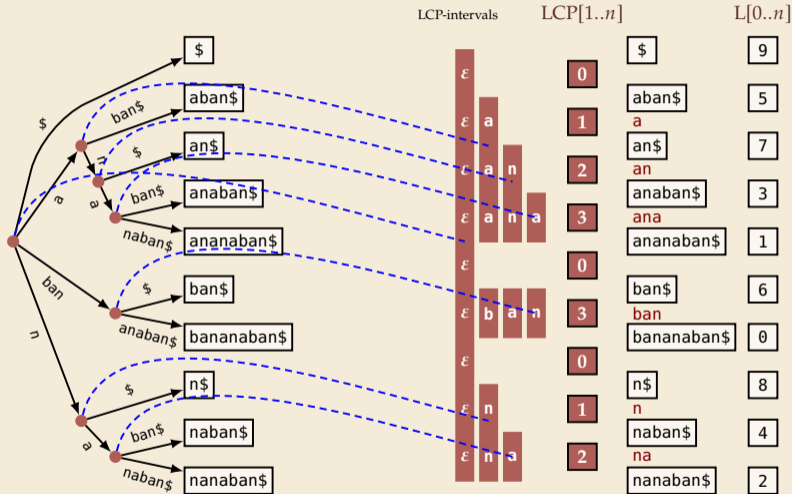
LCP array and internal nodes



LCP array and internal nodes



LCP array and internal nodes



↪ Leaf array $L[0..n]$ plus LCP array $LCP[1..n]$ encode full tree!

LCP array construction

- ▶ computing $LCP[1..n]$ naively too expensive
 - ▶ each value could take $\Theta(n)$ time
- 👎 $\Theta(n^2)$ in total

LCP array construction

- ▶ computing $LCP[1..n]$ naively too expensive
 - ▶ each value could take $\Theta(n)$ time
- 👎 $\Theta(n^2)$ in total
- ▶ but: seeing one large (= costly) LCP value \rightsquigarrow can find another large one!
- ▶ Example: $T = \text{Buffalo_buffalo_buffalo_buffalo\$}$

- ▶ first few suffixes in sorted order:

$$T_{L[0]} = \$$$

$$T_{L[1]} = \text{alo_buffalo\$}$$

$$T_{L[2]} = \text{alo_buffalo_buffalo\$}$$

$$\text{alo_buffalo_buffalo} \rightsquigarrow LCP[3] = 19$$

$$T_{L[3]} = \text{alo_buffalo_buffalo_buffalo\$}$$

LCP array construction

- ▶ computing $LCP[1..n]$ naively too expensive
 - ▶ each value could take $\Theta(n)$ time
- 👎 $\Theta(n^2)$ in total
- ▶ but: seeing one large (= costly) LCP value \rightsquigarrow can find another large one!
- ▶ Example: $T = \text{Buffalo_buffalo_buffalo_buffalo\$}$

- ▶ first few suffixes in sorted order:

$T_{L[0]} = \$$

$T_{L[1]} = \text{alo_buffalo\$}$

$T_{L[2]} = \text{alo_buffalo_buffalo\$}$

alo_buffalo_buffalo $\rightsquigarrow LCP[3] = 19$

$T_{L[3]} = \text{alo_buffalo_buffalo_buffalo\$}$

- \rightsquigarrow **Removing first character** from $T_{L[2]}$ and $T_{L[3]}$ gives two new suffixes:

$T_{L[?]} = \text{lo_buffalo_buffalo\$}$

lo_buffalo_buffalo $\rightsquigarrow LCP[?] = 18$

$T_{L[?]} = \text{lo_buffalo_buffalo_buffalo\$}$

↑
unclear where...

LCP array construction

- ▶ computing $LCP[1..n]$ naively too expensive
 - ▶ each value could take $\Theta(n)$ time
- 👎 $\Theta(n^2)$ in total
- ▶ but: seeing one large (= costly) LCP value \rightsquigarrow can find another large one!
- ▶ Example: $T = \text{Buffalo_buffalo_buffalo_buffalo\$}$

- ▶ first few suffixes in sorted order:

$T_{L[0]} = \$$

$T_{L[1]} = \text{alo_buffalo\$}$

$T_{L[2]} = \text{alo_buffalo_buffalo\$}$

alo_buffalo_buffalo $\rightsquigarrow LCP[3] = 19$

$T_{L[3]} = \text{alo_buffalo_buffalo_buffalo\$}$

- \rightsquigarrow **Removing first character** from $T_{L[2]}$ and $T_{L[3]}$ gives two new suffixes:

$T_{L[?]} = \text{lo_buffalo_buffalo\$}$

lo_buffalo_buffalo $\rightsquigarrow LCP[?] = 18$

$T_{L[?]} = \text{lo_buffalo_buffalo_buffalo\$}$

↑
unclear where...



Shortened suffixes might *not* be *adjacent* in sorted order!

\rightsquigarrow no LCP entry for them!

Kasai's algorithm – Example

- ▶ Kasai et al. used above observation systematically
- ▶ Key idea: *compute* LCP values in *text order*
- ▶ Dropping first character of adjacent suffixes might not lead to *adjacent* shorter suffixes, but LCP entry can only be *longer*.

i	$R[i]$	T_i	r	$L[r]$	$T_{L[r]}$	LCP[r]
0	6 th	bananaban\$	0	9	\$	–
1	4 th	ananaban\$	1	5	aban\$	
2	9 th	nanaban\$	2	7	an\$	
3	3 th	anaban\$	3	3	anaban\$	
4	8 th	naban\$	4	1	ananaban\$	
5	1 th	aban\$	5	6	ban\$	
6	5 th	ban\$	6	0	bananaban\$	
7	2 th	an\$	7	8	n\$	
8	7 th	n\$	8	4	naban\$	
9	0 th	\$	9	2	nanaban\$	

Kasai's algorithm – Example

- ▶ Kasai et al. used above observation systematically
- ▶ Key idea: *compute* LCP values in *text order*
- ▶ Dropping first character of adjacent suffixes might not lead to *adjacent* shorter suffixes, but LCP entry can only be *longer*.

i	$R[i]$	T_i	r	$L[r]$	$T_{L[r]}$	LCP[r]
→ 0	6 th	bananaban\$	0	9	\$	–
1	4 th	ananaban\$	1	5	aban\$	
2	9 th	nanaban\$	2	7	an\$	
3	3 th	anaban\$	3	3	anaban\$	
4	8 th	naban\$	4	1	ananaban\$	
5	1 th	aban\$	5	6	ban\$	
6	5 th	ban\$	6	0	bananaban\$	
7	2 th	an\$	7	8	n\$	
8	7 th	n\$	8	4	naban\$	
9	0 th	\$	9	2	nanaban\$	

Kasai's algorithm – Example

- ▶ Kasai et al. used above observation systematically
- ▶ Key idea: *compute* LCP values in *text order*
- ▶ Dropping first character of adjacent suffixes might not lead to *adjacent* shorter suffixes, but LCP entry can only be *longer*.

i	$R[i]$	T_i	r	$L[r]$	$T_{L[r]}$	LCP[r]
→ 0	6 th	bananaban\$	0	9	\$	-
1	4 th	ananaban\$	1	5	aban\$	
2	9 th	nanaban\$	2	7	an\$	
3	3 th	anaban\$	3	3	anaban\$	
4	8 th	naban\$	4	1	ananaban\$	
5	1 th	aban\$	5	6	ban\$	
6	5 th	ban\$	→ 6	0	bananaban\$	
7	2 th	an\$	7	8	n\$	
8	7 th	n\$	8	4	naban\$	
9	0 th	\$	9	2	nanaban\$	

Kasai's algorithm – Example

- ▶ Kasai et al. used above observation systematically
- ▶ Key idea: *compute* LCP values in *text order*
- ▶ Dropping first character of adjacent suffixes might not lead to *adjacent* shorter suffixes, but LCP entry can only be *longer*.

i	$R[i]$	T_i	r	$L[r]$	$T_{L[r]}$	$LCP[r]$
→ 0	6 th	bananaban\$	0	9	\$	-
1	4 th	ananaban\$	1	5	aban\$	
2	9 th	nanaban\$	2	7	an\$	
3	3 th	anaban\$	3	3	anaban\$	
4	8 th	naban\$	4	1	ananaban\$	
5	1 th	aban\$	5	6	ban\$	
6	5 th	ban\$	→ 6	0	bananaban\$	
7	2 th	an\$	7	8	n\$	
8	7 th	n\$	8	4	naban\$	
9	0 th	\$	9	2	nanaban\$	

Kasai's algorithm – Example

- ▶ Kasai et al. used above observation systematically
- ▶ Key idea: *compute* LCP values in *text order*
- ▶ Dropping first character of adjacent suffixes might not lead to *adjacent* shorter suffixes, but LCP entry can only be *longer*.

i	$R[i]$	T_i	r	$L[r]$	$T_{L[r]}$	$LCP[r]$
→ 0	6 th	bananaban\$	0	9	\$	-
1	4 th	ananaban\$	1	5	aban\$	
2	9 th	nanaban\$	2	7	an\$	
3	3 th	anaban\$	3	3	anaban\$	
4	8 th	naban\$	4	1	ananaban\$	
5	1 th	aban\$	5	6	<u>ban</u> \$	
6	5 th	ban\$	→ 6	0	<u>ban</u> ananaban\$	
7	2 th	an\$	7	8	n\$	
8	7 th	n\$	8	4	naban\$	
9	0 th	\$	9	2	nanaban\$	

Kasai's algorithm – Example

- ▶ Kasai et al. used above observation systematically
- ▶ Key idea: *compute* LCP values in *text order*
- ▶ Dropping first character of adjacent suffixes might not lead to *adjacent* shorter suffixes, but LCP entry can only be *longer*.

i	$R[i]$	T_i	r	$L[r]$	$T_{L[r]}$	$LCP[r]$
→ 0	6 th	bananaban\$	0	9	\$	-
1	4 th	ananaban\$	1	5	aban\$	
2	9 th	nanaban\$	2	7	an\$	
3	3 th	anaban\$	3	3	anaban\$	
4	8 th	naban\$	4	1	ananaban\$	
5	1 th	aban\$	5	6	<u>ban</u> \$	
6	5 th	ban\$	→ 6	0	<u>ban</u> anaban\$	
7	2 th	an\$	7	8	n\$	
8	7 th	n\$	8	4	naban\$	
9	0 th	\$	9	2	nanaban\$	

Kasai's algorithm – Example

- ▶ Kasai et al. used above observation systematically
- ▶ Key idea: *compute* LCP values in *text order*
- ▶ Dropping first character of adjacent suffixes might not lead to *adjacent* shorter suffixes, but LCP entry can only be *longer*.

i	$R[i]$	T_i	r	$L[r]$	$T_{L[r]}$	$LCP[r]$
→ 0	6 th	bananaban\$	0	9	\$	-
1	4 th	ananaban\$	1	5	aban\$	
2	9 th	nanaban\$	2	7	an\$	
3	3 th	anaban\$	3	3	anaban\$	
4	8 th	naban\$	4	1	ananaban\$	
5	1 th	aban\$	5	6	<u>ban</u> \$	
6	5 th	ban\$	→ 6	0	<u>ban</u> anaban\$	3
7	2 th	an\$	7	8	n\$	
8	7 th	n\$	8	4	naban\$	
9	0 th	\$	9	2	nanaban\$	

Kasai's algorithm – Example

- ▶ Kasai et al. used above observation systematically
- ▶ Key idea: *compute* LCP values in *text order*
- ▶ Dropping first character of adjacent suffixes might not lead to *adjacent* shorter suffixes, but LCP entry can only be *longer*.

i	$R[i]$	T_i	r	$L[r]$	$T_{L[r]}$	LCP[r]
0	6 th	bananaban\$	0	9	\$	–
1	4 th	ananaban\$	1	5	aban\$	
2	9 th	nanaban\$	2	7	an\$	
3	3 th	anaban\$	3	3	anaban\$	
4	8 th	naban\$	4	1	ananaban\$	
5	1 th	aban\$	5	6	ban \$	
6	5 th	ban\$	6	0	ban anaban\$	3
7	2 th	an\$	7	8	n\$	
8	7 th	n\$	8	4	naban\$	
9	0 th	\$	9	2	nanaban\$	

Kasai's algorithm – Example

- ▶ Kasai et al. used above observation systematically
- ▶ Key idea: *compute* LCP values in *text order*
- ▶ Dropping first character of adjacent suffixes might not lead to *adjacent* shorter suffixes, but LCP entry can only be *longer*.

i	$R[i]$	T_i	r	$L[r]$	$T_{L[r]}$	LCP[r]
0	6 th	bananaban\$	0	9	\$	–
→ 1	4 th	ananaban\$	1	5	aban\$	
2	9 th	nanaban\$	2	7	an\$	
3	3 th	anaban\$	3	3	anaban\$	
4	8 th	naban\$	4	1	ananaban\$	
5	1 th	aban\$	5	6	ban \$	
6	5 th	ban\$	6	0	ban anaban\$	3
7	2 th	an\$	7	8	n\$	
8	7 th	n\$	8	4	naban\$	
9	0 th	\$	9	2	nanaban\$	

Kasai's algorithm – Example

- ▶ Kasai et al. used above observation systematically
- ▶ Key idea: *compute* LCP values in *text order*
- ▶ Dropping first character of adjacent suffixes might not lead to *adjacent* shorter suffixes, but LCP entry can only be *longer*.

i	$R[i]$	T_i	r	$L[r]$	$T_{L[r]}$	$LCP[r]$
0	6 th	bananaban\$	0	9	\$	-
→ 1	4 th	ananaban\$	1	5	aban\$	
2	9 th	nanaban\$	2	7	an\$	
3	3 th	anaban\$	3	3	anaban\$	
4	8 th	naban\$	→ 4	1	ananaban\$	
5	1 th	aban\$	5	6	<u>ba</u> n\$	
6	5 th	ban\$	6	0	<u>ba</u> nanaban\$	3
7	2 th	an\$	7	8	n\$	
8	7 th	n\$	8	4	naban\$	
9	0 th	\$	9	2	nanaban\$	

Kasai's algorithm – Example

- ▶ Kasai et al. used above observation systematically
- ▶ Key idea: *compute* LCP values in *text order*
- ▶ Dropping first character of adjacent suffixes might not lead to *adjacent* shorter suffixes, but LCP entry can only be *longer*.

i	$R[i]$	T_i	r	$L[r]$	$T_{L[r]}$	$LCP[r]$
0	6 th	bananaban\$	0	9	\$	-
→ 1	4 th	ananaban\$	1	5	aban\$	
2	9 th	nanaban\$	2	7	an\$	
3	3 th	anaban\$	3	3	anaban\$	
4	8 th	naban\$	→ 4	1	anaban\$	
5	1 th	aban\$	5	6	ban\$	
6	5 th	ban\$	6	0	bananaban\$	3
7	2 th	an\$	7	8	n\$	
8	7 th	n\$	8	4	naban\$	
9	0 th	\$	9	2	nanaban\$	

Kasai's algorithm – Example

- ▶ Kasai et al. used above observation systematically
- ▶ Key idea: *compute* LCP values in *text order*
- ▶ Dropping first character of adjacent suffixes might not lead to *adjacent* shorter suffixes, but LCP entry can only be *longer*.

i	$R[i]$	T_i	r	$L[r]$	$T_{L[r]}$	LCP[r]
0	6 th	bananaban\$	0	9	\$	-
→ 1	4 th	ananaban\$	1	5	aban\$	
2	9 th	nanaban\$	2	7	an\$	
3	3 th	anaban\$	3	3	anaban\$	
4	8 th	naban\$	→ 4	1	ananaban\$	
5	1 th	aban\$	5	6	ban\$	
6	5 th	ban\$	6	0	bananaban\$	3
7	2 th	an\$	7	8	n\$	
8	7 th	n\$	8	4	naban\$	
9	0 th	\$	9	2	nanaban\$	

Kasai's algorithm – Example

- ▶ Kasai et al. used above observation systematically
- ▶ Key idea: *compute* LCP values in *text order*
- ▶ Dropping first character of adjacent suffixes might not lead to *adjacent* shorter suffixes, but LCP entry can only be *longer*.

i	$R[i]$	T_i	r	$L[r]$	$T_{L[r]}$	LCP[r]
0	6 th	bananaban\$	0	9	\$	-
→ 1	4 th	ananaban\$	1	5	aban\$	
2	9 th	nanaban\$	2	7	an\$	
3	3 th	anaban\$	3	3	anaban\$	
4	8 th	naban\$	→ 4	1	ananaban\$	3
5	1 th	aban\$	5	6	ban\$	
6	5 th	ban\$	6	0	bananaban\$	3
7	2 th	an\$	7	8	n\$	
8	7 th	n\$	8	4	naban\$	
9	0 th	\$	9	2	nanaban\$	

Kasai's algorithm – Example

- ▶ Kasai et al. used above observation systematically
- ▶ Key idea: *compute* LCP values in *text order*
- ▶ Dropping first character of adjacent suffixes might not lead to *adjacent* shorter suffixes, but LCP entry can only be *longer*.

i	$R[i]$	T_i	r	$L[r]$	$T_{L[r]}$	$LCP[r]$
0	6 th	bananaban\$	0	9	\$	–
1	4 th	ananaban\$	1	5	aban\$	
2	9 th	nanaban\$	2	7	an\$	
3	3 th	anaban\$	3	3	an aban\$	
4	8 th	naban\$	4	1	an anaban\$	3
5	1 th	aban\$	5	6	ban\$	
6	5 th	ban\$	6	0	bananaban\$	3
7	2 th	an\$	7	8	n\$	
8	7 th	n\$	8	4	naban\$	
9	0 th	\$	9	2	nanaban\$	

Kasai's algorithm – Example

- ▶ Kasai et al. used above observation systematically
- ▶ Key idea: *compute* LCP values in *text order*
- ▶ Dropping first character of adjacent suffixes might not lead to *adjacent* shorter suffixes, but LCP entry can only be *longer*.

i	$R[i]$	T_i	r	$L[r]$	$T_{L[r]}$	$LCP[r]$
0	6 th	bananaban\$	0	9	\$	–
1	4 th	ananaban\$	1	5	aban\$	
→ 2	9 th	nanaban\$	2	7	an\$	
3	3 th	anaban\$	3	3	anaban\$	
4	8 th	naban\$	4	1	anaban\$	3
5	1 th	aban\$	5	6	ban\$	
6	5 th	ban\$	6	0	bananaban\$	3
7	2 th	an\$	7	8	n\$	
8	7 th	n\$	8	4	naban\$	
9	0 th	\$	9	2	nanaban\$	

Kasai's algorithm – Example

- ▶ Kasai et al. used above observation systematically
- ▶ Key idea: *compute* LCP values in *text order*
- ▶ Dropping first character of adjacent suffixes might not lead to *adjacent* shorter suffixes, but LCP entry can only be *longer*.

i	$R[i]$	T_i	r	$L[r]$	$T_{L[r]}$	LCP[r]
0	6 th	bananaban\$	0	9	\$	–
1	4 th	ananaban\$	1	5	aban\$	
→ 2	9 th	nanaban\$	2	7	an\$	
3	3 th	anaban\$	3	3	an <u>a</u> ban\$	
4	8 th	naban\$	4	1	an <u>a</u> nanaban\$	3
5	1 th	aban\$	5	6	ban\$	
6	5 th	ban\$	6	0	bananaban\$	3
7	2 th	an\$	7	8	n\$	
8	7 th	n\$	8	4	naban\$	
9	0 th	\$	→ 9	2	nanaban\$	

Kasai's algorithm – Example

- ▶ Kasai et al. used above observation systematically
- ▶ Key idea: *compute* LCP values in *text order*
- ▶ Dropping first character of adjacent suffixes might not lead to *adjacent* shorter suffixes, but LCP entry can only be *longer*.

i	$R[i]$	T_i	r	$L[r]$	$T_{L[r]}$	$LCP[r]$
0	6 th	bananaban\$	0	9	\$	–
1	4 th	ananaban\$	1	5	aban\$	
→ 2	9 th	nanaban\$	2	7	an\$	
3	3 th	anaban\$	3	3	anaban\$	
4	8 th	naban\$	4	1	anaban\$	3
5	1 th	aban\$	5	6	ban\$	
6	5 th	ban\$	6	0	bananaban\$	3
7	2 th	an\$	7	8	n\$	
8	7 th	n\$	8	4	naban\$	
9	0 th	\$	→ 9	2	nanaban\$	

Kasai's algorithm – Example

- ▶ Kasai et al. used above observation systematically
- ▶ Key idea: *compute* LCP values in *text order*
- ▶ Dropping first character of adjacent suffixes might not lead to *adjacent* shorter suffixes, but LCP entry can only be *longer*.

i	$R[i]$	T_i	r	$L[r]$	$T_{L[r]}$	LCP[r]
0	6 th	bananaban\$	0	9	\$	–
1	4 th	ananaban\$	1	5	aban\$	
→ 2	9 th	nanaban\$	2	7	an\$	
3	3 th	anaban\$	3	3	an <u>a</u> ban\$	
4	8 th	naban\$	4	1	an <u>a</u> nanaban\$	3
5	1 th	aban\$	5	6	ban\$	
6	5 th	ban\$	6	0	bananaban\$	3
7	2 th	an\$	7	8	n\$	
8	7 th	n\$	8	4	<u>n</u> aban\$	
9	0 th	\$	→ 9	2	<u>n</u> anaban\$	2

Kasai's algorithm – Example

- ▶ Kasai et al. used above observation systematically
- ▶ Key idea: *compute* LCP values in *text order*
- ▶ Dropping first character of adjacent suffixes might not lead to *adjacent* shorter suffixes, but LCP entry can only be *longer*.

i	$R[i]$	T_i	r	$L[r]$	$T_{L[r]}$	LCP[r]
0	6 th	bananaban\$	0	9	\$	–
1	4 th	ananaban\$	1	5	aban\$	
2	9 th	nanaban\$	2	7	an\$	
3	3 th	anaban\$	3	3	anaban\$	
4	8 th	naban\$	4	1	ananaban\$	3
5	1 th	aban\$	5	6	ban\$	
6	5 th	ban\$	6	0	bananaban\$	3
7	2 th	an\$	7	8	n\$	
8	7 th	n\$	8	4	naban\$	
9	0 th	\$	9	2	nanaban\$	2

Kasai's algorithm – Example

- ▶ Kasai et al. used above observation systematically
- ▶ Key idea: *compute* LCP values in *text order*
- ▶ Dropping first character of adjacent suffixes might not lead to *adjacent* shorter suffixes, but LCP entry can only be *longer*.

i	$R[i]$	T_i	r	$L[r]$	$T_{L[r]}$	LCP[r]
0	6 th	bananaban\$	0	9	\$	–
1	4 th	ananaban\$	1	5	aban\$	
2	9 th	nanaban\$	2	7	an\$	
→ 3	3 th	anaban\$	3	3	anaban\$	
4	8 th	naban\$	4	1	ananaban\$	3
5	1 th	aban\$	5	6	ban\$	
6	5 th	ban\$	6	0	bananaban\$	3
7	2 th	an\$	7	8	n\$	
8	7 th	n\$	8	4	naban\$	
9	0 th	\$	9	2	nanaban\$	2

Kasai's algorithm – Example

- ▶ Kasai et al. used above observation systematically
- ▶ Key idea: *compute* LCP values in *text order*
- ▶ Dropping first character of adjacent suffixes might not lead to *adjacent* shorter suffixes, but LCP entry can only be *longer*.

i	$R[i]$	T_i	r	$L[r]$	$T_{L[r]}$	$LCP[r]$
0	6 th	bananaban\$	0	9	\$	–
1	4 th	ananaban\$	1	5	aban\$	
2	9 th	nanaban\$	2	7	an\$	
→ 3	3 th	anaban\$	→ 3	3	anaban\$	
4	8 th	naban\$	4	1	ananaban\$	3
5	1 th	aban\$	5	6	ban\$	
6	5 th	ban\$	6	0	bananaban\$	3
7	2 th	an\$	7	8	n\$	
8	7 th	n\$	8	4	naban\$	
9	0 th	\$	9	2	nanaban\$	2

Kasai's algorithm – Example

- ▶ Kasai et al. used above observation systematically
- ▶ Key idea: *compute* LCP values in *text order*
- ▶ Dropping first character of adjacent suffixes might not lead to *adjacent* shorter suffixes, but LCP entry can only be *longer*.

i	$R[i]$	T_i	r	$L[r]$	$T_{L[r]}$	LCP[r]
0	6 th	bananaban\$	0	9	\$	–
1	4 th	ananaban\$	1	5	aban\$	
2	9 th	nanaban\$	2	7	an\$	
→ 3	3 th	anaban\$	→ 3	3	anaban\$	
4	8 th	naban\$	4	1	ananaban\$	3
5	1 th	aban\$	5	6	ban\$	
6	5 th	ban\$	6	0	bananaban\$	3
7	2 th	an\$	7	8	n\$	
8	7 th	n\$	8	4	naban\$	
9	0 th	\$	9	2	nanaban\$	2

Kasai's algorithm – Example

- ▶ Kasai et al. used above observation systematically
- ▶ Key idea: *compute* LCP values in *text order*
- ▶ Dropping first character of adjacent suffixes might not lead to *adjacent* shorter suffixes, but LCP entry can only be *longer*.

i	$R[i]$	T_i	r	$L[r]$	$T_{L[r]}$	$LCP[r]$
0	6 th	bananaban\$	0	9	\$	–
1	4 th	ananaban\$	1	5	aban\$	
2	9 th	nanaban\$	2	7	an\$	
→ 3	3 th	anaban\$	→ 3	3	anaban\$	
4	8 th	naban\$	4	1	ananaban\$	3
5	1 th	aban\$	5	6	ban\$	
6	5 th	ban\$	6	0	bananaban\$	3
7	2 th	an\$	7	8	n\$	
8	7 th	n\$	8	4	naban\$	
9	0 th	\$	9	2	nanaban\$	2

Kasai's algorithm – Example

- ▶ Kasai et al. used above observation systematically
- ▶ Key idea: *compute* LCP values in *text order*
- ▶ Dropping first character of adjacent suffixes might not lead to *adjacent* shorter suffixes, but LCP entry can only be *longer*.

i	$R[i]$	T_i	r	$L[r]$	$T_{L[r]}$	$LCP[r]$
0	6 th	bananaban\$	0	9	\$	–
1	4 th	ananaban\$	1	5	aban\$	
2	9 th	nanaban\$	2	7	an\$	
→ 3	3 th	anaban\$	→ 3	3	anaban\$	2
4	8 th	naban\$	4	1	ananaban\$	3
5	1 th	aban\$	5	6	ban\$	
6	5 th	ban\$	6	0	bananaban\$	3
7	2 th	an\$	7	8	n\$	
8	7 th	n\$	8	4	naban\$	
9	0 th	\$	9	2	nanaban\$	2

Kasai's algorithm – Example

- ▶ Kasai et al. used above observation systematically
- ▶ Key idea: *compute* LCP values in *text order*
- ▶ Dropping first character of adjacent suffixes might not lead to *adjacent* shorter suffixes, but LCP entry can only be *longer*.

i	$R[i]$	T_i	r	$L[r]$	$T_{L[r]}$	LCP[r]
0	6 th	bananaban\$	0	9	\$	–
1	4 th	ananaban\$	1	5	aban\$	
2	9 th	nanaban\$	2	7	an\$	
3	3 th	anaban\$	3	3	anaban\$	2
4	8 th	naban\$	4	1	ananaban\$	3
5	1 th	aban\$	5	6	ban\$	
6	5 th	ban\$	6	0	bananaban\$	3
7	2 th	an\$	7	8	n\$	
8	7 th	n\$	8	4	naban\$	
9	0 th	\$	9	2	nanaban\$	2

Kasai's algorithm – Example

- ▶ Kasai et al. used above observation systematically
- ▶ Key idea: *compute* LCP values in *text order*
- ▶ Dropping first character of adjacent suffixes might not lead to *adjacent* shorter suffixes, but LCP entry can only be *longer*.

i	$R[i]$	T_i	r	$L[r]$	$T_{L[r]}$	$LCP[r]$
0	6 th	bananaban\$	0	9	\$	–
1	4 th	ananaban\$	1	5	aban\$	
2	9 th	nanaban\$	2	7	an\$	
3	3 th	anaban\$	3	3	anaban\$	2
→ 4	8 th	naban\$	4	1	ananaban\$	3
5	1 th	aban\$	5	6	ban\$	
6	5 th	ban\$	6	0	bananaban\$	3
7	2 th	an\$	7	8	n\$	
8	7 th	n\$	→ 8	4	naban\$	
9	0 th	\$	9	2	nanaban\$	2

Kasai's algorithm – Example

- ▶ Kasai et al. used above observation systematically
- ▶ Key idea: *compute* LCP values in *text order*
- ▶ Dropping first character of adjacent suffixes might not lead to *adjacent* shorter suffixes, but LCP entry can only be *longer*.

i	$R[i]$	T_i	r	$L[r]$	$T_{L[r]}$	LCP[r]
0	6 th	bananaban\$	0	9	\$	–
1	4 th	ananaban\$	1	5	aban\$	
2	9 th	nanaban\$	2	7	an\$	
3	3 th	anaban\$	3	3	anaban\$	2
→ 4	8 th	naban\$	4	1	ananaban\$	3
5	1 th	aban\$	5	6	ban\$	
6	5 th	ban\$	6	0	bananaban\$	3
7	2 th	an\$	7	8	n\$	
8	7 th	n\$	→ 8	4	naban\$	1
9	0 th	\$	9	2	nanaban\$	2

Kasai's algorithm – Example

- ▶ Kasai et al. used above observation systematically
- ▶ Key idea: *compute* LCP values in *text order*
- ▶ Dropping first character of adjacent suffixes might not lead to *adjacent* shorter suffixes, but LCP entry can only be *longer*.

i	$R[i]$	T_i	r	$L[r]$	$T_{L[r]}$	$LCP[r]$
0	6 th	bananaban\$	0	9	\$	–
1	4 th	ananaban\$	→ 1	5	aban\$	
2	9 th	nanaban\$	2	7	an\$	
3	3 th	anaban\$	3	3	anaban\$	2
4	8 th	naban\$	4	1	ananaban\$	3
→ 5	1 th	aban\$	5	6	ban\$	
6	5 th	ban\$	6	0	bananaban\$	3
7	2 th	an\$	7	8	n\$	
8	7 th	n\$	8	4	naban\$	1
9	0 th	\$	9	2	nanaban\$	2

Kasai's algorithm – Example

- ▶ Kasai et al. used above observation systematically
- ▶ Key idea: *compute* LCP values in *text order*
- ▶ Dropping first character of adjacent suffixes might not lead to *adjacent* shorter suffixes, but LCP entry can only be *longer*.

i	$R[i]$	T_i	r	$L[r]$	$T_{L[r]}$	$LCP[r]$
0	6 th	bananaban\$	0	9	\$	–
1	4 th	ananaban\$	→ 1	5	aban\$	0
2	9 th	nanaban\$	2	7	an\$	
3	3 th	anaban\$	3	3	anaban\$	2
4	8 th	naban\$	4	1	ananaban\$	3
→ 5	1 th	aban\$	5	6	ban\$	
6	5 th	ban\$	6	0	bananaban\$	3
7	2 th	an\$	7	8	n\$	
8	7 th	n\$	8	4	naban\$	1
9	0 th	\$	9	2	nanaban\$	2

Kasai's algorithm – Example

- ▶ Kasai et al. used above observation systematically
- ▶ Key idea: *compute* LCP values in *text order*
- ▶ Dropping first character of adjacent suffixes might not lead to *adjacent* shorter suffixes, but LCP entry can only be *longer*.

i	$R[i]$	T_i	r	$L[r]$	$T_{L[r]}$	$LCP[r]$
0	6 th	bananaban\$	0	9	\$	–
1	4 th	ananaban\$	1	5	aban\$	0
2	9 th	nanaban\$	2	7	an\$	
3	3 th	anaban\$	3	3	anaban\$	2
4	8 th	naban\$	4	1	ananaban\$	3
5	1 th	aban\$	5	6	ban\$	
→ 6	5 th	ban\$	6	0	bananaban\$	3
7	2 th	an\$	7	8	n\$	
8	7 th	n\$	8	4	naban\$	1
9	0 th	\$	9	2	nanaban\$	2

Kasai's algorithm – Example

- ▶ Kasai et al. used above observation systematically
- ▶ Key idea: *compute* LCP values in *text order*
- ▶ Dropping first character of adjacent suffixes might not lead to *adjacent* shorter suffixes, but LCP entry can only be *longer*.

i	$R[i]$	T_i	r	$L[r]$	$T_{L[r]}$	LCP[r]
0	6 th	bananaban\$	0	9	\$	–
1	4 th	ananaban\$	1	5	aban\$	0
2	9 th	nanaban\$	2	7	an\$	
3	3 th	anaban\$	3	3	anaban\$	2
4	8 th	naban\$	4	1	ananaban\$	3
5	1 th	aban\$	5	6	ban\$	0
6	5 th	ban\$	6	0	bananaban\$	3
7	2 th	an\$	7	8	n\$	
8	7 th	n\$	8	4	naban\$	1
9	0 th	\$	9	2	nanaban\$	2

Kasai's algorithm – Example

- ▶ Kasai et al. used above observation systematically
- ▶ Key idea: *compute* LCP values in *text order*
- ▶ Dropping first character of adjacent suffixes might not lead to *adjacent* shorter suffixes, but LCP entry can only be *longer*.

i	$R[i]$	T_i	r	$L[r]$	$T_{L[r]}$	$LCP[r]$
0	6 th	bananaban\$	0	9	\$	–
1	4 th	ananaban\$	1	5	aban\$	0
2	9 th	nanaban\$	2	7	an\$	
3	3 th	anaban\$	3	3	anaban\$	2
4	8 th	naban\$	4	1	ananaban\$	3
5	1 th	aban\$	5	6	ban\$	0
6	5 th	ban\$	6	0	bananaban\$	3
7	2 th	an\$	7	8	n\$	
8	7 th	n\$	8	4	naban\$	1
9	0 th	\$	9	2	nanaban\$	2

Kasai's algorithm – Example

- ▶ Kasai et al. used above observation systematically
- ▶ Key idea: *compute* LCP values in *text order*
- ▶ Dropping first character of adjacent suffixes might not lead to *adjacent* shorter suffixes, but LCP entry can only be *longer*.

i	$R[i]$	T_i	r	$L[r]$	$T_{L[r]}$	$LCP[r]$
0	6 th	bananaban\$	0	9	\$	–
1	4 th	ananaban\$	1	5	aban\$	0
2	9 th	nanaban\$	2	7	an\$	1
3	3 th	anaban\$	3	3	anaban\$	2
4	8 th	naban\$	4	1	ananaban\$	3
5	1 th	aban\$	5	6	ban\$	0
6	5 th	ban\$	6	0	bananaban\$	3
7	2 th	an\$	7	8	n\$	
8	7 th	n\$	8	4	naban\$	1
9	0 th	\$	9	2	nanaban\$	2

Kasai's algorithm – Example

- ▶ Kasai et al. used above observation systematically
- ▶ Key idea: *compute* LCP values in *text order*
- ▶ Dropping first character of adjacent suffixes might not lead to *adjacent* shorter suffixes, but LCP entry can only be *longer*.

i	$R[i]$	T_i	r	$L[r]$	$T_{L[r]}$	LCP[r]
0	6 th	bananaban\$	0	9	\$	–
1	4 th	ananaban\$	1	5	aban\$	0
2	9 th	nanaban\$	2	7	an\$	1
3	3 th	anaban\$	3	3	anaban\$	2
4	8 th	naban\$	4	1	ananaban\$	3
5	1 th	aban\$	5	6	ban\$	0
6	5 th	ban\$	6	0	bananaban\$	3
7	2 th	an\$	7	8	n\$	
8	7 th	n\$	8	4	naban\$	1
9	0 th	\$	9	2	nanaban\$	2

Kasai's algorithm – Example

- ▶ Kasai et al. used above observation systematically
- ▶ Key idea: *compute* LCP values in *text order*
- ▶ Dropping first character of adjacent suffixes might not lead to *adjacent* shorter suffixes, but LCP entry can only be *longer*.

i	$R[i]$	T_i	r	$L[r]$	$T_{L[r]}$	$LCP[r]$
0	6 th	bananaban\$	0	9	\$	–
1	4 th	ananaban\$	1	5	aban\$	0
2	9 th	nanaban\$	2	7	an\$	1
3	3 th	anaban\$	3	3	anaban\$	2
4	8 th	naban\$	4	1	ananaban\$	3
5	1 th	aban\$	5	6	ban\$	0
6	5 th	ban\$	6	0	bananaban\$	3
7	2 th	an\$	7	8	n\$	0
8	7 th	n\$	8	4	naban\$	1
9	0 th	\$	9	2	nanaban\$	2

Kasai's algorithm – Code

```
1 procedure computeLCP( $T[0..n]$ ,  $L[0..n]$ ,  $R[0..n]$ )
2   // Assume  $T[n] = \$$ ,  $L$  and  $R$  are suffix array and inverse
3    $\ell := 0$ 
4   for  $i := 0, \dots, n-1$  // looks at  $T_i$ 
5      $r := R[i]$ 
6     // compute LCP[ $r$ ]; note that  $r > 0$  since  $R[n] = 0$ 
7      $i_{-1} := L[r-1]$ 
8     while  $T[i+\ell] == T[i_{-1}+\ell]$  do
9        $\ell := \ell + 1$ 
10    LCP[ $r$ ] :=  $\ell$ 
11     $\ell := \max\{\ell - 1, 0\}$ 
12  return LCP[ $1..n$ ]
```

- ▶ remember length ℓ of induced common prefix
- ▶ use L to get start index of suffixes

Kasai's algorithm – Code

```
1 procedure computeLCP( $T[0..n]$ ,  $L[0..n]$ ,  $R[0..n]$ )
2   // Assume  $T[n] = \$$ ,  $L$  and  $R$  are suffix array and inverse
3    $\ell := 0$ 
4   for  $i := 0, \dots, n - 1$ 
5      $r := R[i]$ 
6     // compute  $LCP[r]$ ; note that  $r > 0$  since  $R[n] = 0$ 
7      $i_{-1} := L[r - 1]$ 
8     while  $T[i + \ell] == T[i_{-1} + \ell]$  do
9        $\ell := \ell + 1$ 
10     $LCP[r] := \ell$ 
11     $\ell := \max\{\ell - 1, 0\}$ 
12  return  $LCP[1..n]$ 
```

- ▶ remember length ℓ of induced common prefix
- ▶ use L to get start index of suffixes

Analysis:

- ▶ dominant operation: character comparisons
- ▶ separately count those with outcomes “=” resp. “≠”
- ▶ each ≠ ends iteration of for-loop
 $\rightsquigarrow \leq n$ cmps
- ▶ each = implies increment of ℓ , but $\ell \leq n$ and decremented $\leq n$ times
 $\rightsquigarrow \leq 2n$ cmps
- $\rightsquigarrow \Theta(n)$ overall time

Conclusion

▶ (*Enhanced*) *Suffix Arrays* are the modern version of suffix trees

👎 can be harder to reason about

👍 can support same algorithms as suffix trees

👍 but use much less space

👍 simpler linear-time construction