

13

Text Indexing – Searching entire genomes

2 February 2026

Prof. Dr. Sebastian Wild

Learning Outcomes

Unit 13: *Text Indexing*

1. Know and understand methods for text indexing: *inverted indexes, suffix trees, (enhanced) suffix arrays*
2. Know and understand *generalized suffix trees*
3. Know properties, in particular *performance characteristics*, and limitations of the above data structures.
4. Design (simple) *algorithms based on suffix trees*.
5. Understand *construction algorithms* for suffix arrays.

13 Text Indexing

- 13.1 Inverted Indexes and Tries
- 13.2 Suffix Trees
- 13.3 Applications
- 13.4 Generalized Suffix Trees
- 13.5 Suffix Arrays
- 13.6 Linear-Time Suffix Sorting: Inducing Order
- 13.7 Linear-Time Suffix Sorting: The DC3 Algorithm

13.1 Inverted Indexes and Tries

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 indexes

≈ same as “indices”

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

Inverted indexes

≈ same as “indices”

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

Inverted index:

- ▶ collect all words in T
 - ▶ can be as simple as splitting T at whitespaces
 - ▶ actual implementations typically support *stemming* of words
goes \rightarrow go, cats \rightarrow cat
language specific!
- ▶ store mapping from words to a list of occurrences \rightsquigarrow *how?*

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/cs566*

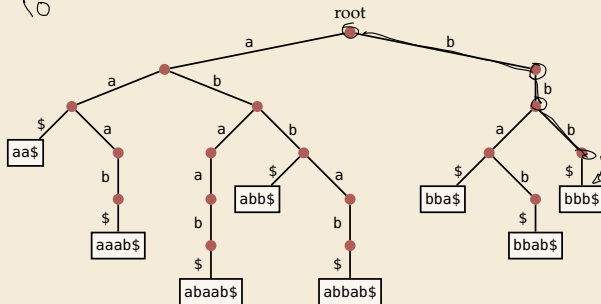
Tries

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

some character $\notin \Sigma$

▶ Example:

{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 $|Q| = q$ (with $q \leq m$).

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/cs566

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 $|Q| = q$ (with $q \leq m$).

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

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

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

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

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

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



→ sli.do/cs566

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 $\Theta(m)$

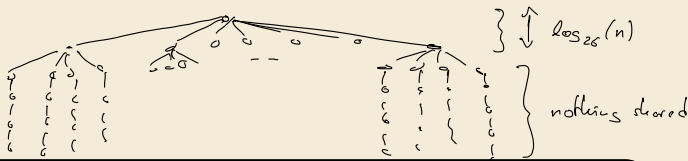
C $\Theta(n \cdot m)$

F $\Theta(m \log n)$



→ sli.do/cs566

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 ~~$\Theta(n + m)$~~ \equiv total # chars

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

C $\Theta(n \cdot m)$ ✓

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

$$n(m - \log n) + n$$

node in trie needs array of child pointers
 $\approx \Theta(5)$ space per node

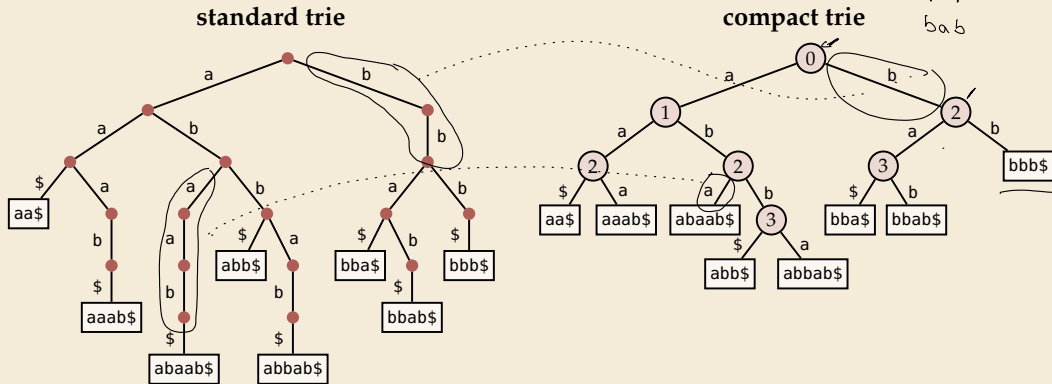


→ sli.do/cs566

Compact tries

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

=1 child



↪ searching slightly trickier, but same time complexity as in trie

- ▶ all nodes ≥ 2 children \rightsquigarrow $\#nodes \leq \#leaves = \#strings \rightsquigarrow$ linear space in $\#strings$

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

```
ACAAGATGCCATTGTCCCCGGCCTCCTGCTGCTGCTCTCCGGGGCCACGGCCACCGCTGCCCTGCCCTGGAGGGTGGCCCCACGGC  
CGAGACAGCGAGCATATGCAGGAAGCGGCAGGAATAAGGAAAAGCAGCCTCCTGACTTTCCTCGCTTGGTGGTTTGAGTGGACCTCCAGGC  
CAGTGCCGGGCCCCCTCATAGGAGAGGAAGCTCGGGAGGTGGCCAGGCGGCAGGAAGGCGACCCCCCAGCAATCCGCGCGCCGGGACAGAA  
TGCCCTGCAGGAACCTTCTTCTGGAAGACCTTCTCCTCCTGCAAATAAAACCTCACCCATGAATGCTCACGCAAGTTTAATTACAGACCTGAA
```

- ▶ binary streams

```
00000010101001111010111000001111100011111011111001101101000011100010011011110000010001101010  
011011000011010110100000001000000011101011000001000011110101110110010001100101101110111111  
110001010001011001010000001110101010011000000001101100001100111110000101 0101011101111000011  
1010111001001010101010100000111110100110000001111001101010000000100100100000101100011000110111
```

~> need new ideas

13.2 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

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

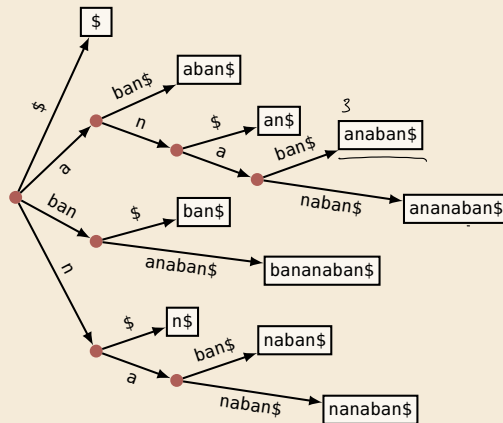
Example:

$T = \text{bananaban\$}$

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

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 copy of actual string)

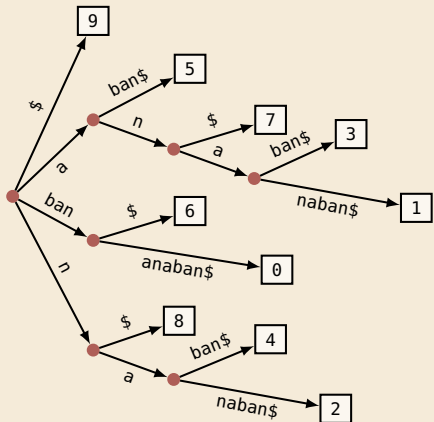
Example:

$T = \text{bananaban\$}$

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
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)

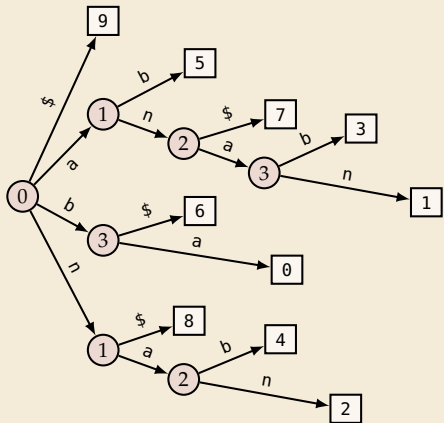
Example:

$T = \text{bananaban\$}$

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

	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!

- ▶ algorithms are a bit tricky to understand \rightarrow EAA Book
- ▶ 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?

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/cs566

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.~~ $n+1$



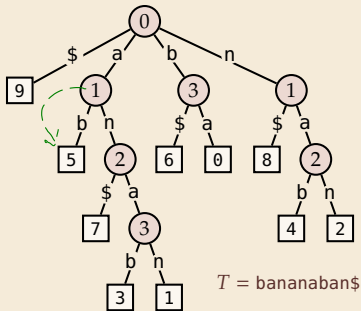
→ sli.do/cs566

13.3 Applications

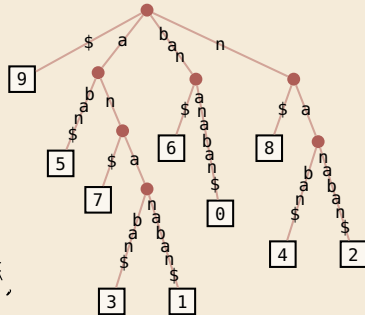
Applications of suffix trees

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

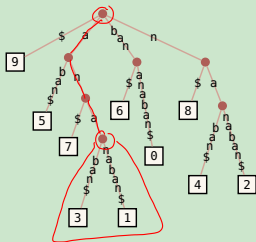
Recall: \mathcal{T} stored like this:



but think about this:




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




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.

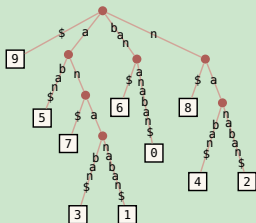
I P starts at index 3. 

J P starts at index 4.

K P starts at index 7.



→ sli.do/cs566



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 occurs three times in T .~~
 - 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/cs566

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

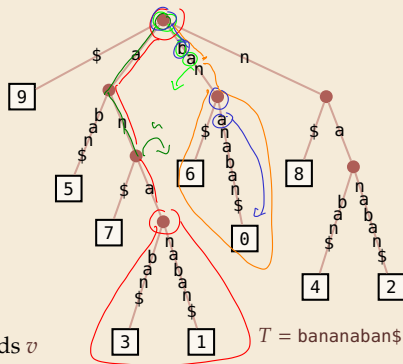
3. **we run out of tree**

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:

► $P = \underline{\text{ann}}$

► $P = \underline{\text{baa}}$

► $P = \underline{\text{ana}}$

► $P = \underline{\text{ba}}$

► $P = \underline{\text{briar}}$

Application 2: Longest repeated substring

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

e.g. for compression \rightsquigarrow Unit 7



How can we efficiently check *all possible substrings*?

► **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*

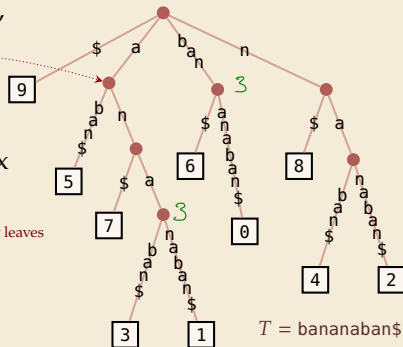


- $\rightsquigarrow \exists$ internal node with path label 'a'

here single edge, can be longer path

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

actually: adjacent leaves



► **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



Repeated substrings = shared paths in *suffix tree*



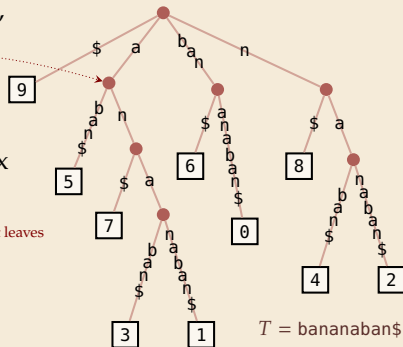
- $\rightsquigarrow \exists$ internal node with path label 'a'

^ here single edge, can be longer path

actually: adjacent leaves

1. Compute *string depth* (=length of path label) of nodes

- Both can be done in depth-first traversal $\rightsquigarrow \Theta(n)$ time



13.4 Generalized Suffix Trees

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

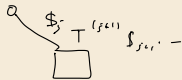
Generalized suffix trees

- ▶ 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 $T^{(j)} \dots$ but doesn't seem to help
- \rightsquigarrow need a *single/joint* suffix tree for *several* texts

Enter: *generalized suffix tree*

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

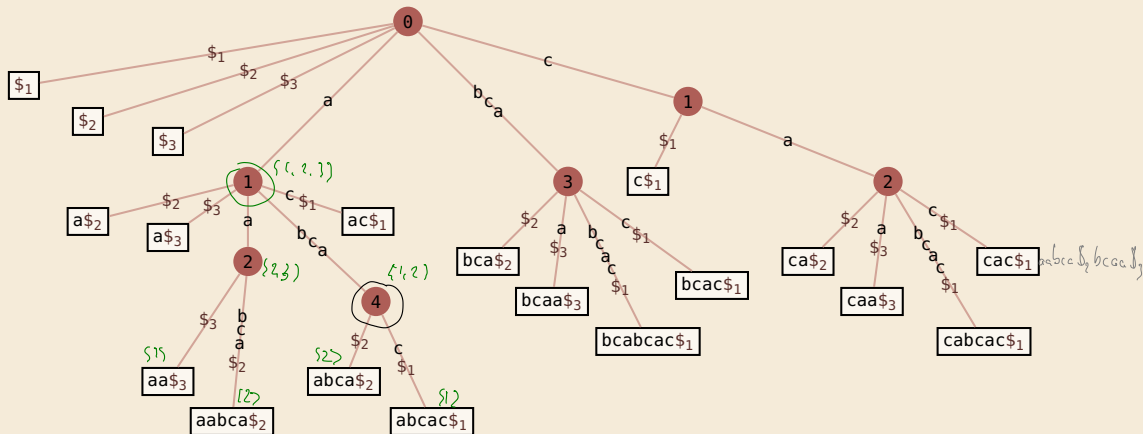
\rightsquigarrow $\$_j$ -edges always leads to leaves $\rightsquigarrow \exists$ leaf $\underline{(j, i)}$ for each suffix $T_i^{(j)} = T^{(j)}[i..n_j]$



Generalized Suffix Tree – Example

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

$T = \text{bcabcac}\$, \text{aabca}\$, \text{bcaa}\$$



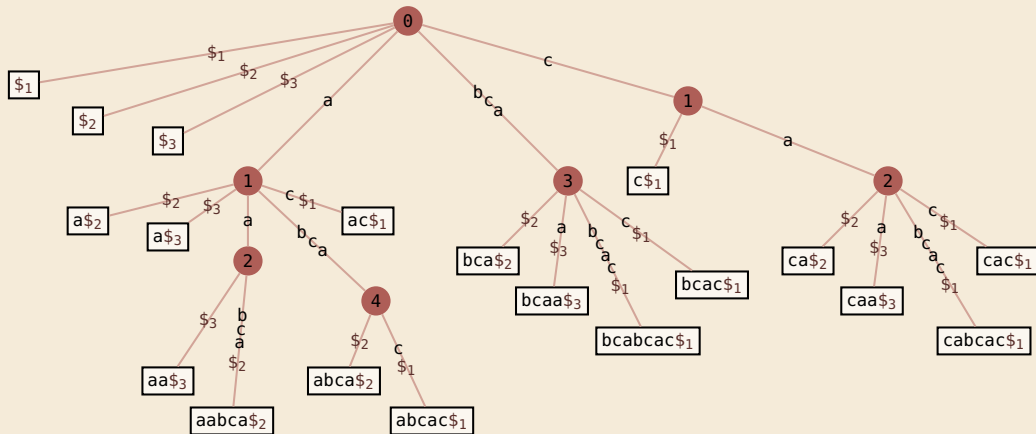
Application 3: Longest common substring

- ▶ 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 $\{1, \dots, k\}$.
- ▶ Each step takes time $\Theta(n)$ for $n = n_1 + \dots + n_k$ the total length of all texts. (for constant k)

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

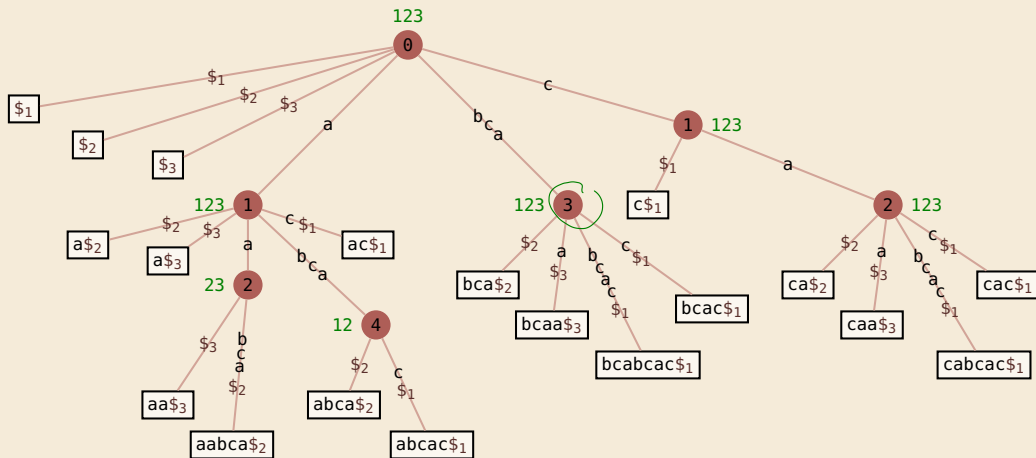
Longest common substring – Example

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



Longest common substring – Example

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



Further applications

Many more applications, in particular for problems on biological sequences

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

- ▶ key ingredient: longest common extensions

If you want to see more, come to *Algorithms of Bioinformatics* 😊

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



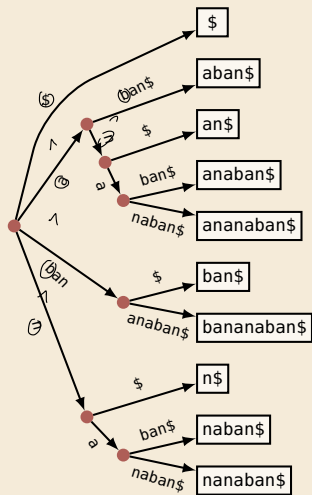
many pointers in tree incur large space overhead



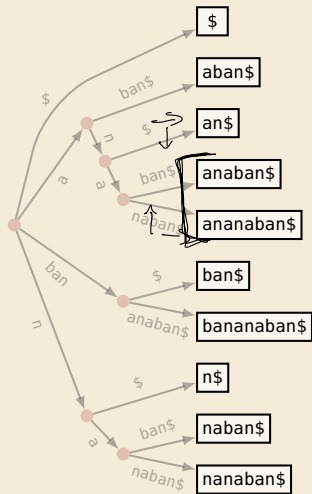
13.5 Suffix Arrays

Putting suffix trees on a diet

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

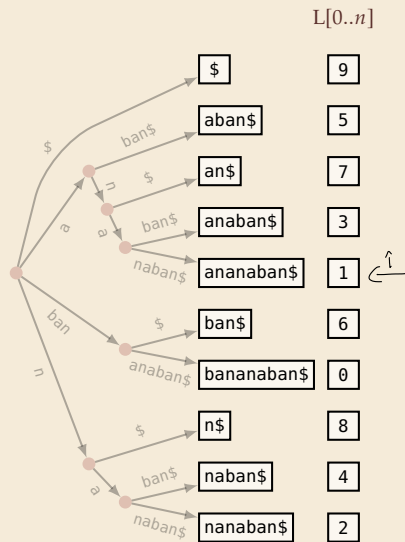


Putting suffix trees on a diet



- ▶ **Observation:** order of leaves in suffix tree
= suffixes lexicographically *sorted*
- ▶ Idea: only store list of leaves $L[0..n]$
- ▶ 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:** $P = \text{ana}$

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
 $\leq (\lg n + 2) \cdot m$ character comparisons

Clicker Question



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

- A** L can be very easily computed from B and T
- B** B can be very easily computed from L and T
- C** Both A and B
- D** Neither A nor B



→ sli.do/cs566

Clicker Question



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

- ☐ **A** ~~L can be very easily computed from B and T~~
- ☒ **B** B can be very easily computed from L and T ✓
- ☐ **C** ~~Both A and B~~
- ☐ **D** ~~Neither A nor B~~



→ sli.do/cs566

Digression: Recall BWT

Burrows-Wheeler Transform

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

$S = \text{alf_eats_alfalfa\$}$

$B = \text{asff\$f_e_lllaaata}$

[illegible]

sort

\$alf_eats_alfalfa
_alfalfa\$alf_eats
_eats_alfalfa\$alf
a\$alf_eats_alfalfa
alf_eats_alfalfa
alfalfa\$alf_eats_alf
alfalfa\$alf_eats_alf
eats_alfalfa\$alf_e
eats_alfalfa\$alf_e
f_eats_alfalfa\$alf_e
fa\$alf_eats_alfalfa
falfa\$alf_eats_alfalfa
lf_eats_alfalfa\$alf_e
lfa\$alf_eats_alfalfa
lfalfa\$alf_eats_alf
s_alfalfa\$alf_eats
ts_alfalfa\$alf_eats

BWT



Digression: Computing the BWT

How can we compute the BWT of a text efficiently?

Digression: Computing the BWT

How can we compute the BWT of a text efficiently?

- ▶ cyclic shifts $S \hat{=}$ suffixes of S
 - ▶ comparing cyclic shifts stops at first \$
 - ▶ for comparisons, anything after \$ irrelevant!
- ▶ BWT is essentially suffix sorting!
 - ▶ $B[i] = S[L[i] - 1]$
 - ▶ where $L[i] = 0, B[i] = \$$

↪ Can compute B in $O(n)$ time from L

```

alf_eats_alfalfa$
lf_eats_alfalfa$
f_eats_alfalfa$al
_eats_alfalfa$alf
eats_alfalfa$alf_
ats_alfalfa$alf_e
ts_alfalfa$alf_ea
s_alfalfa$alf_eat
_alfalfa$alf_eats
alfalfa$alf_eats_
lfalfa$alf_eats_a
falfa$alf_eats_al
alfa$alf_eats_alf
lfa$alf_eats_alfa
fa$alf_eats_alfal
a$alf_eats_alfalf
$alf_eats_alfalfa
  
```

	r		$L[r]$
0	\$alf_eats_alfalfa	a	16
1	_alfalfa\$alf_eats	s	8
2	_eats_alfalfa\$alf	f	3
3	a\$alf_eats_alfalf	f	15
4	alf_eats_alfalfa\$	\$	0
5	alfa\$alf_eats_alf	f	12
6	alfalfa\$alf_eats_	_	9
7	ats_alfalfa\$alf_e	e	5
8	eats_alfalfa\$alf_	_	4
9	f_eats_alfalfa\$al	l	2
10	fa\$alf_eats_alfal	l	14
11	falfa\$alf_eats_al	l	11
12	lf_eats_alfalfa\$a	a	1
13	lfa\$alf_eats_alfa	a	13
14	lfalfa\$alf_eats_a	a	10
15	s_alfalfa\$alf_eat	t	7
16	ts_alfalfa\$alf_ea	a	6

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” 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 *Ex 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



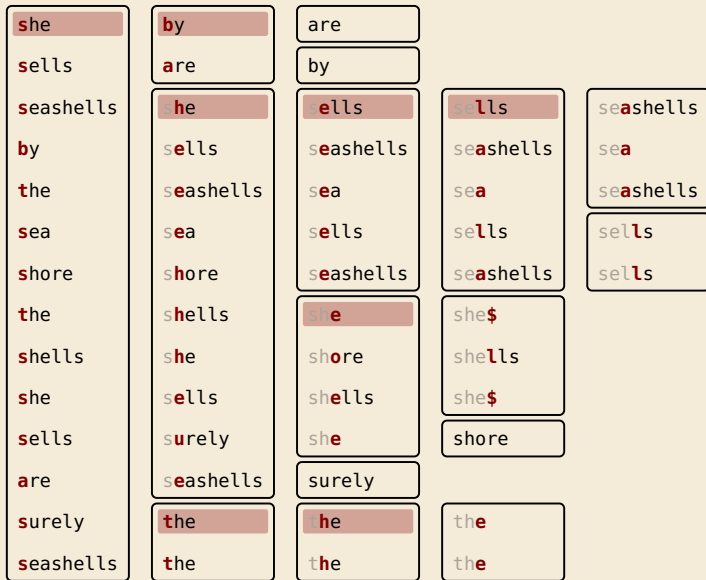
Fat-pivot radix quicksort – Example

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	

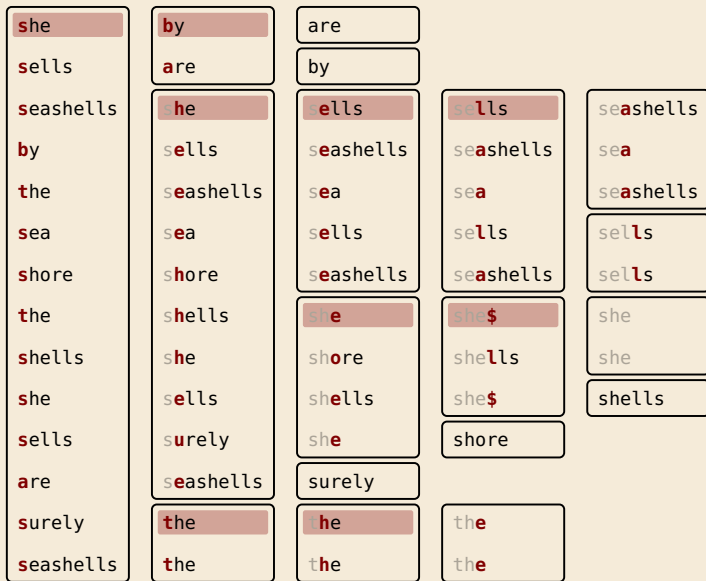
Fat-pivot radix quicksort – Example

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



Fat-pivot radix quicksort – Example



Fat-pivot radix quicksort – Example



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
 - ▶ duplicates only reduce # 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
- ▶ duplicates only reduce # 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**

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

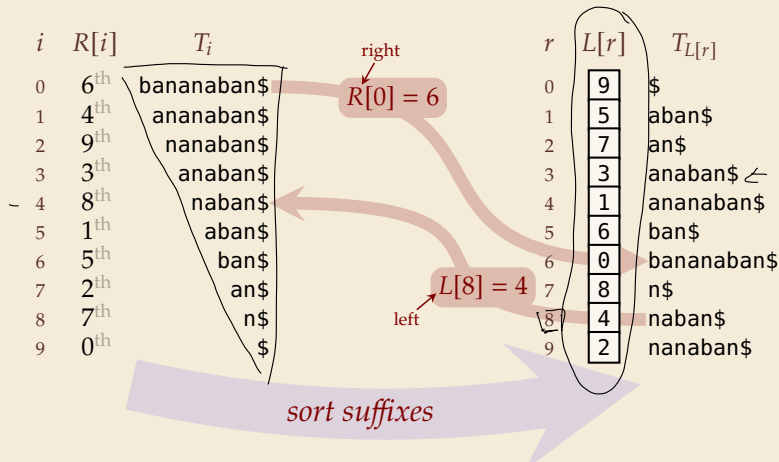
but we can do $O(n)$ time worst case!

13.6 Linear-Time Suffix Sorting: Inducing Order

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*



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/cs566

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/cs566

Linear-time suffix sorting

DC3 / Skew algorithm

not a multiple of 3



1. Compute rank array $R_{1,2}$ for suffixes T_i starting at $i \not\equiv 0 \pmod{3}$ 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

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
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 $R_{1,2}$.
↪ rank array R for entire input

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

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

Linear-time suffix sorting

DC3 / Skew algorithm

not a multiple of 3



1. Compute rank array $R_{1,2}$ for suffixes T_i starting at $i \not\equiv 0 \pmod{3}$ *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

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

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

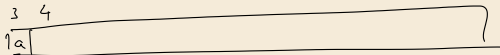
\rightsquigarrow Can use whichever is more convenient.

DC3 / Skew algorithm – Step 2: Inducing ranks

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

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

- **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 .

- ▶ 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 $R_{1,2}[1]$ and $R_{1,2}[4]$!

~> 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 sbanasman\$\$\$
 T_{12} nanasman\$\$\$
 T_{15} asman\$\$\$
 T_{18} an\$\$\$
 T_{21} \$\$

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

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

sman\$\$\$ =

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

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

 $R_1, 2$ (known)

radix sort

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 sbanasman\$\$\$
 T_{12} nanasman\$\$\$
 T_{15} asman\$\$\$
 T_{18} an\$\$\$
 T_{21} \$\$

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

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

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

radix sort

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

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

$R_{1,2}$ (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 sbanasman\$\$\$
 T_{12} nanasman\$\$\$
 T_{15} asman\$\$\$
 T_{18} an\$\$\$
 T_{21} \$\$

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

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

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

radix sort

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

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

$R_{1,2}$ (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

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 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 $R_{1,2}$

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

↪ Compare T_{16} to T_{12}

$T_{16} = \text{sman}$$$
 $= \text{sman}$$$
 $= sT_{17}$
 $T_{12} = \text{nanasman}$$$
 $= \text{aananasman}$$$
 $= 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 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}$$$
 $= \text{asman}$$$ can't compare T_{16}
 $= aT_{16}$ and T_{12} either!
 $T_{11} = \text{ananasman}$$$
 $= \text{ananasman}$$$
 $= aT_{12}$$$$$

↪ Compare T_{16} to T_{12}

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

$\rightsquigarrow O(n)$ time for merge

Compare T_{15} to T_{11}

Idea: try same trick as before

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

\rightsquigarrow Compare T_{16} to T_{12}

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

continue 15:16

13.7 Linear-Time Suffix Sorting: The DC3 Algorithm

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* abc

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

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

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

$T = \text{bananaban} \text{ $$$ }$
↪ $T^\square = \begin{array}{l} \text{ban} \text{ ana} \text{ ban} \text{ $$$} \\ \text{ana} \text{ ban} \text{ $$$} \\ \text{ban} \text{ $$$} \\ \text{$$$} \end{array}$

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!

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

$$T' = T[1..n)^{\square} \boxed{\$ \$ \$} T[2..n)^{\square} \boxed{\$ \$ \$}$$

► $T = \text{hannahbansbananasman\$}$

DC3 / Skew algorithm – Step 3. Example

$$T' = T[1..n)^{\square} \boxed{\$ \$ \$} T[2..n)^{\square} \boxed{\$ \$ \$}$$

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

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

DC3 / Skew algorithm – Step 3. Example

$$T' = T[1..n)^{\square} \boxed{\$ \$ \$} T[2..n)^{\square} \boxed{\$ \$ \$}$$

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

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

► Occurring triples:

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

DC3 / Skew algorithm – Step 3. Example

$$T' = T[1..n)^{\square} \text{\texttt{\$ \$ \$ \$}} T[2..n)^{\square} \text{\texttt{\$ \$ \$ \$}}$$

► $T = \text{hannahbansbananasman\$}$ $T_2 = \text{nnahbansbananasman\$}$

$T' =$ annahbansbananasman\$\$ \$\$\$ nnahbansbananasman\$\$\$

► Occurring triples:

annahbansbananasman\$\$ \$\$\$ nnahbansb 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

$$T' = T[1..n)^{\square} \boxed{\$ \$ \$} T[2..n)^{\square} \boxed{\$ \$ \$}$$

► $T = \text{hannahbansbananasman\$}$ $T_2 = \text{nnahbansbananasman\$}$

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

► Occurring triples:

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

► Sorted triples with ranks:

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

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

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

Suffix array – Discussion

- 👍 sleek data structure compared to suffix tree
- 👍 simple and fast $O(n \log n)$ construction
- 👍 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



Outlook: Enhanced Suffix Arrays

- ▶ suffix array by itself somewhat less powerful, but can be augmented with the LCP array

longest common prefix

↪ *(Enhanced) Suffix Arrays*

- ▶ the modern version of suffix trees
- ▶ directly simulate suffix tree operations on L and LCP arrays

👎 can be harder to reason about

👍 can support same algorithms as suffix trees

👍 use less space

👍 simpler linear-time construction

↪ Basis for modern *compressed self-indexes* such as *FM index*