

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.

Outline

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 index:

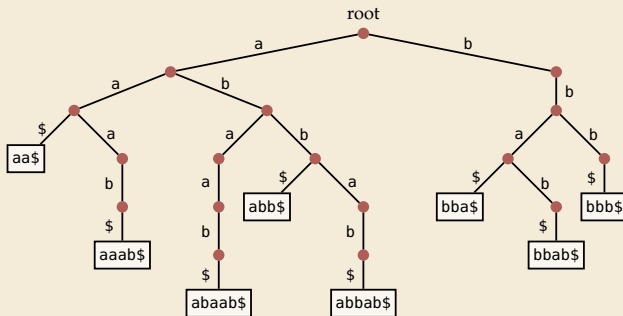
- ▶ collect all words in T
 - ▶ can be as simple as splitting T at whitespaces
 - ▶ actual implementations typically support *stemming* of words
goes \rightarrow go, cats \rightarrow cat
language specific!
- ▶ store mapping from words to a list of occurrences \rightsquigarrow *how?*

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 \$ ✓

- ▶ **Example:**

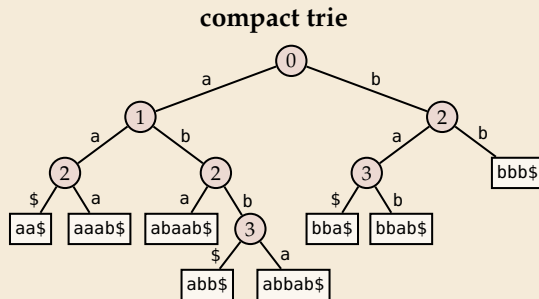
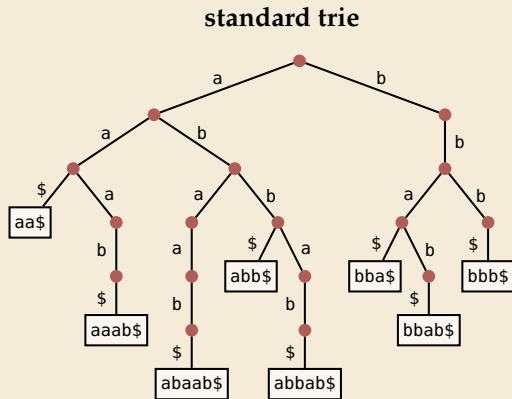
{aa\$, aaab\$, abaab\$, abb\$,
abbab\$, bba\$, bbab\$, bbb\$}



Compact tries

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

=1 child



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

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

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

```
ACAAGATGCCATTGTCCCCGGCCTCCTGCTGCTGCTCTCCGGGGCCACGGCCACCGCTGCCCTGCCCTGGAGGGTGGCCCCACCGGC  
CGAGACAGCGAGCATATGCAGGAAGCGGCAGGAATAAGGAAAAGCAGCCTCCTGACTTTCCTCGCTTGGTGGTTTGAGTGGACCTCCAGGC  
CAGTGCCGGGCCCCCATAGGAGAGGAAGCTCGGGAGGTGGCCAGGCGGCAGGAAGGCGCACCCCCAGCAATCCGCGCGCCGGGACAGAA  
TGCCCTGCAGGAACCTTCTCTGGAAGACCTTCTCTCCTGCAAATAAAACCTCACCCATGAATGCTCACGCAAGTTTAATTACAGACCTGAA
```

- ▶ binary streams

```
00000010101001111010111000001111100011111011111001101101000011100010011011110000010001101010  
0110110000110101101000000010000000011101011000001000011110101110110010001100101101111111  
110001010001011001010000001110101010011000000001101100001100111110000101 0101011101111000011  
10101110010010101010100000111110100110000001111001101010000000100100100000101100011000110111
```

~> 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??

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 $T[n] := \$$)
- ▶ except: in leaves, store *start index* (instead of copy of actual string)

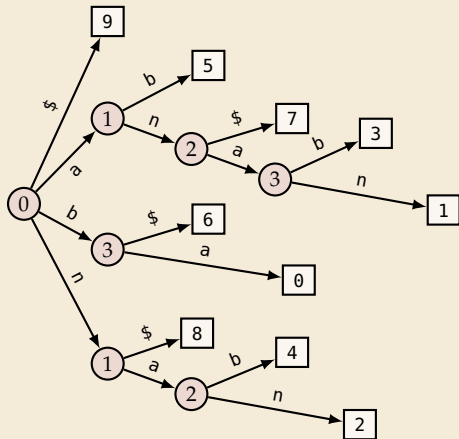
Example:

$T = \text{bananaban\$}$

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	\$

- ▶ 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!



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)!

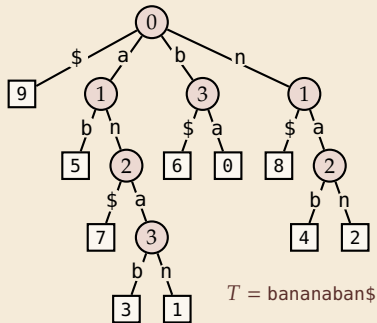
\rightsquigarrow for now, take linear-time construction for granted. What can we do with them?

13.3 Applications

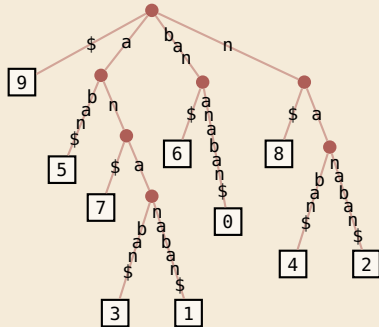
Applications of suffix trees

- 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 \$)

Application 1: Text Indexing / String Matching

► P occurs in $T \iff P$ is a prefix of a suffix of T

► we have all suffixes in T !

↪ (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

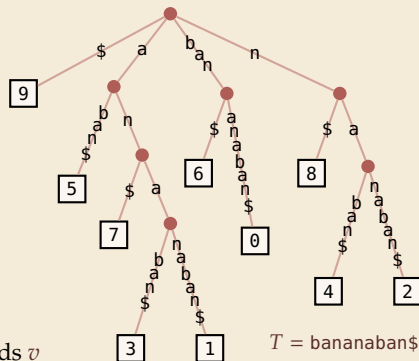
3. we run out of tree

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:

- $P = \text{ann}$
- $P = \text{baa}$
- $P = \text{ana}$
- $P = \text{ba}$
- $P = \text{briar}$

Application 2: Longest repeated substrings

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



How can we efficiently check *all possible substrings*?



Repeated substrings = shared paths in *suffix tree*



e. g. for compression \rightsquigarrow Unit 7

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

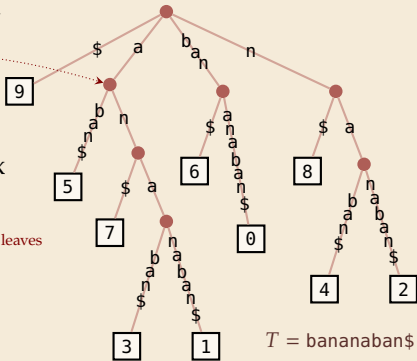
\rightsquigarrow longest repeated substring = longest common prefix (LCP) of two suffixes

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



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
- \rightsquigarrow need a *single/joint* suffix tree for *several* texts

Enter: *generalized suffix tree*

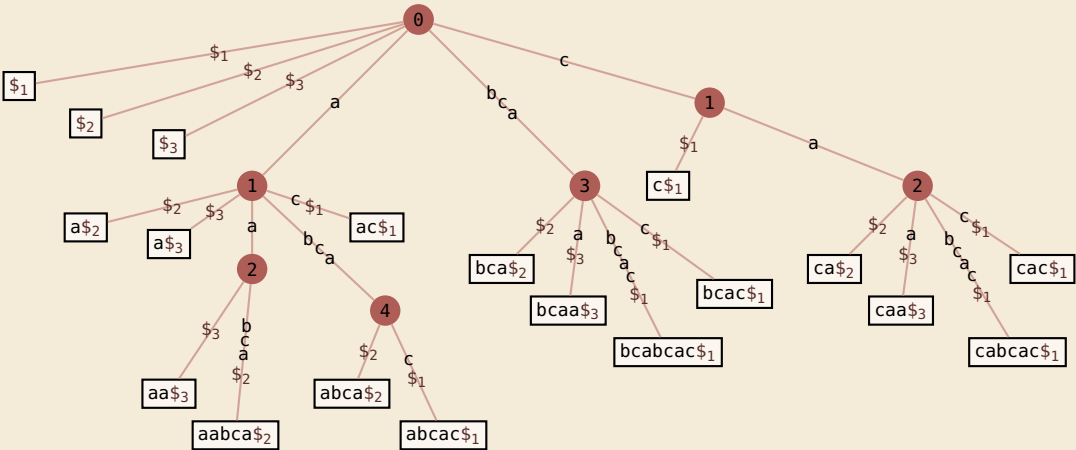
- ▶ Define $T := T^{(1)}\$_1 T^{(2)}\$_2 \dots T^{(k)}\$_k$ for k new end-of-word symbols
- ▶ Construct suffix tree \mathcal{T} for T

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



Generalized Suffix Tree – Example

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



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 $\{1, \dots, k\}$.

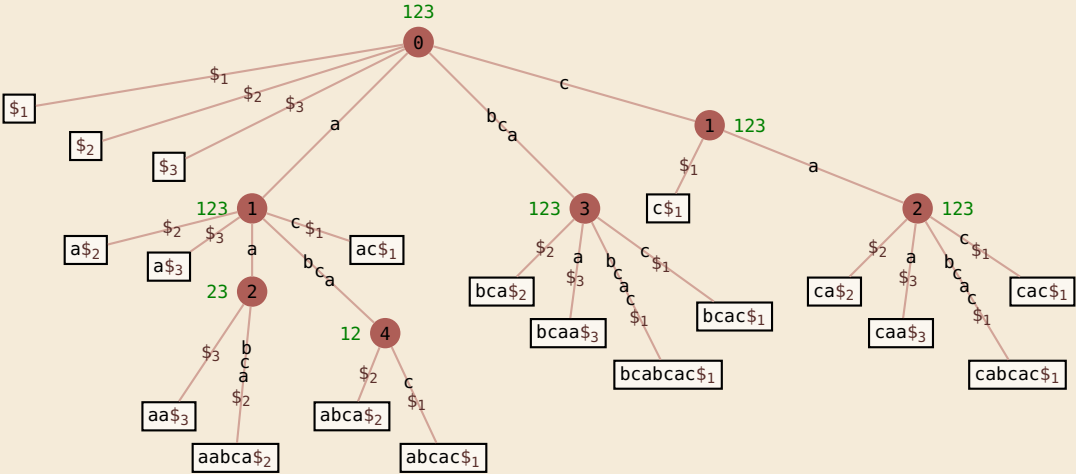
- ▶ Each step takes time $\Theta(n)$ for $n = n_1 + \dots + n_k$ the total length of all texts. (for constant k)

“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

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



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



construction of suffix trees:
linear time, but significant overhead



construction methods fairly complicated



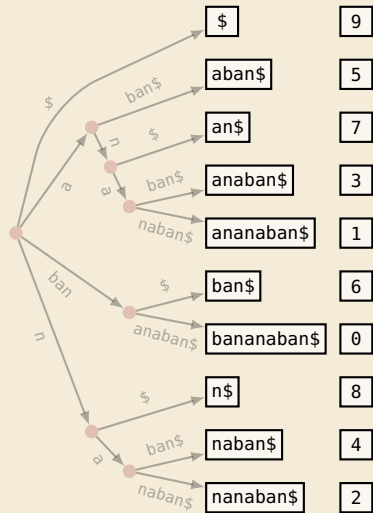
many pointers in tree incur large space overhead



13.5 Suffix Arrays

Putting suffix trees on a diet

$L[0..n]$



► **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

Digression: Recall BWT

Burrows-Wheeler Transform

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

$S = \text{alf_eats_alfalfa\$}$

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

alf_eats_alfalfa\$
lf_eats_alfalfa\$
f_eats_alfalfa\$
_eats_alfalfa\$
eats_alfalfa\$
ats_alfalfa\$
ts_alfalfa\$
s_alfalfa\$
_alfalfa\$
alfalfa\$
lfalfa\$
falffa\$
alfa\$
lfa\$
fa\$
a\$
\$alf_eats_alfalfa

sort

\$alf_eats_alfalf**a**
_alfalfa\$alf_eat**s**
_eats_alfalfa\$alf**f**
a\$alf_eats_alfalf**f**
alf_eats_alfalfa\$**f**
alfalfa\$alf_eats_alf**f**
alfalfa\$alf_eats_**e**
ats_alfalfa\$alf_**t**
eats_alfalfa\$alf_**t**
f_eats_alfalfa\$a**l**
fa\$alf_eats_alfalf**a**
falffa\$alf_eats_alf**a**
lf_eats_alfalfa\$a**a**
lfa\$alf_eats_alf**a**
lfalfa\$alf_eats_**a**
s_alfalfa\$alf_eat**t**
ts_alfalfa\$alf_ea**a**

BWT
↓

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 $L[i] = 0, B[i] = \$$

\rightsquigarrow Can compute B in $O(n)$ time from L

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

(length of) longest common prefix

- ▶ Option 1: Mergesort with LCP values for adjacent elements in runs
- ▶ 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	by	are					
sells	are	by					
seashells	she	sells	sells	seashells	seashells	seashells	sea
by	sells	seashells	seashells	sea	sea\$	seashells	seashells
the	seashells	sea	sea	sells	seashells	sells	seashells
sea	sea	sells	sells	sells	sells	sells	sells
shore	shore	seashells	seashells	sells	sells	sells	sells
the	shells	she	she\$	she	she	shells	
shells	she	shore	shells	shells	shells		
she	sells	shells	she\$				
sells	surely	she	shore				
are	seashells	surely					
surely	the	the	the	the			
seashells	the	the	the	the			

Fat-pivot radix quicksort – 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
- 👎 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

random string



*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**

i	$R[i]$	T_i		r	$L[r]$	$T_{L[r]}$
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\$

right

$R[0] = 6$

left

$L[8] = 4$

sort suffixes

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 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:

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

- ▶ **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} \$\$

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

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}	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}	smans

$R_{1,2}$ (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	$R_0[9] = 7$

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 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 $R_{1,2}$

$\rightsquigarrow O(n)$ 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}$$$$

\rightsquigarrow 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}$$$
 $= \text{aanasman}$$$
 $= aT_{13}$$$$

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!
 - ▶ 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}

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

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

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

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

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

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

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

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

► Occurring triples:

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

► Sorted triples with ranks:

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

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

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

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


↪ (Enhanced) Suffix Arrays longest common prefix

- ▶ 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*