

Applied Databases

Lecture 16

Suffix Array, Burrows-Wheeler Transform

Sebastian Maneth

University of Edinburgh - March 10th, 2016

Outline

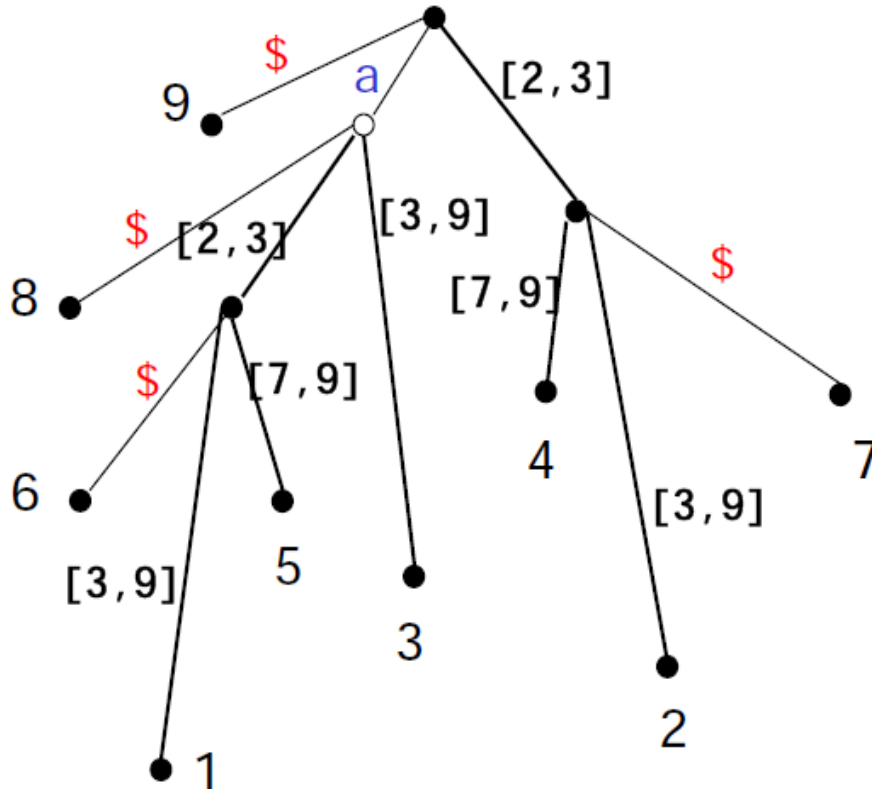
1. Suffix Array
2. Burrows-Wheeler Transform

Online String-Matching

→ how can we do Horspool on Unicode Texts?

Suffix Tree

123456789
 T = abaababa\$



→ add **end marker "\$"**

→ one-to-one correspondence of leaves to suffixes

→ a tree with **m+1** leaves has $\leq 2m+1$ nodes!

Lemma

Size of suffix tree for "**T\$**" is linear in $n=|T|$, i.e., in $O(n)$.

→ search time still $O(|P|)$, as for suffix trie!
 → perfect data structure for our task!

Suffix Tree Construction

Good news:

Suffix tree *can* be constructed in linear time!

Complex construction algorithms

→ Weiner 1973

→ McCreight 1976

→ Ukkonen 1995

→ Farach 1997

Linear time only for *constant-size alphabets!*
Otherwise, $O(n \log n)$

Linear time for **any integer alphabet**,
drawn from a polynomial range

→ again a big breakthrough

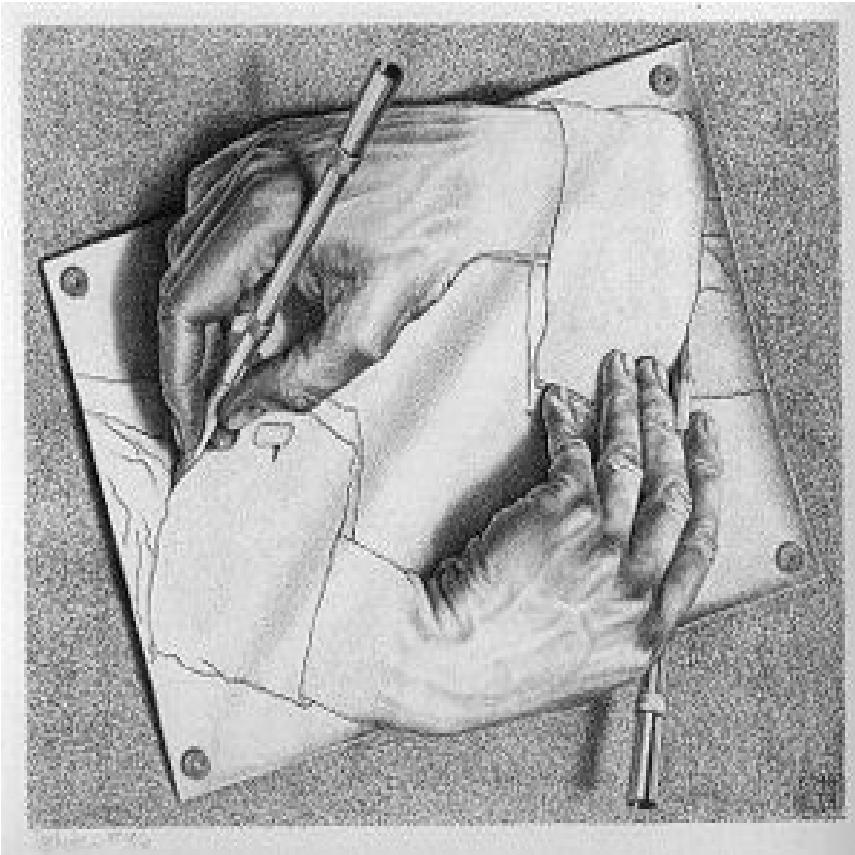
Suffix Tree

Many applications

→ E.g., **infinite-window Lempel-Ziv like compression**

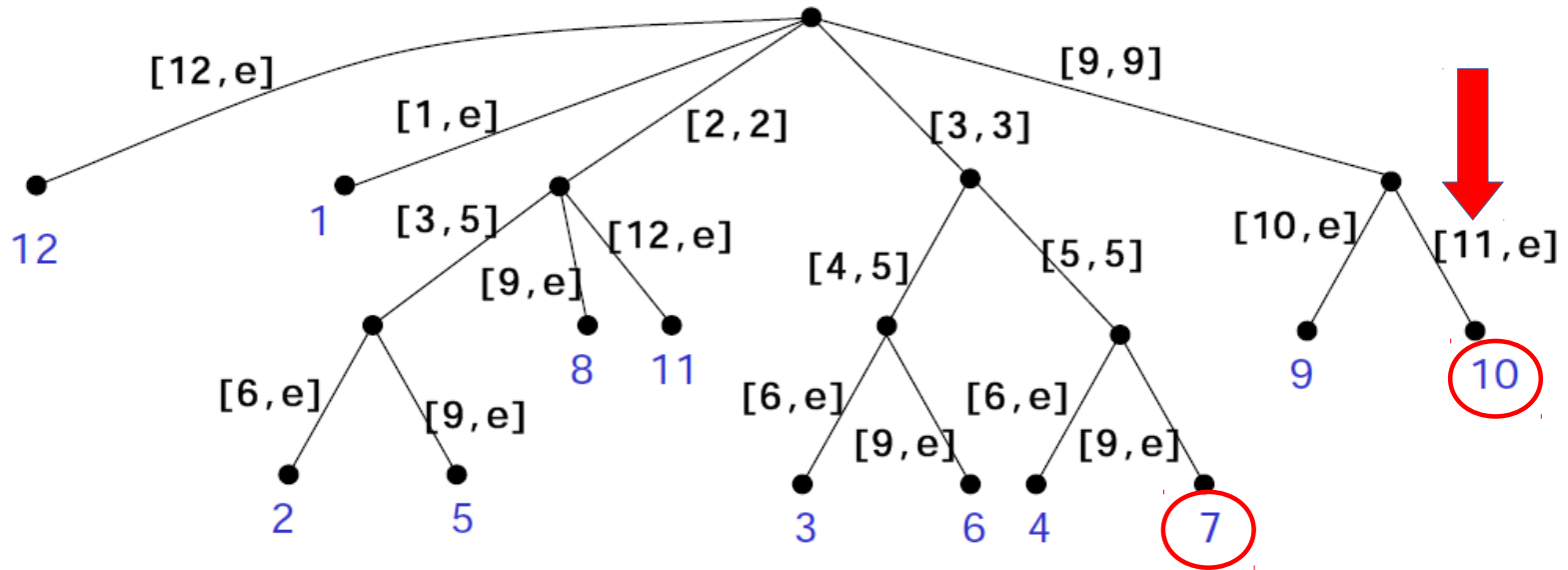
a b a abaa aba baba ab b → a b a (1,4) (1,3) (9,4) (1,2) b

(position, length)



M. C. Escher (1948)

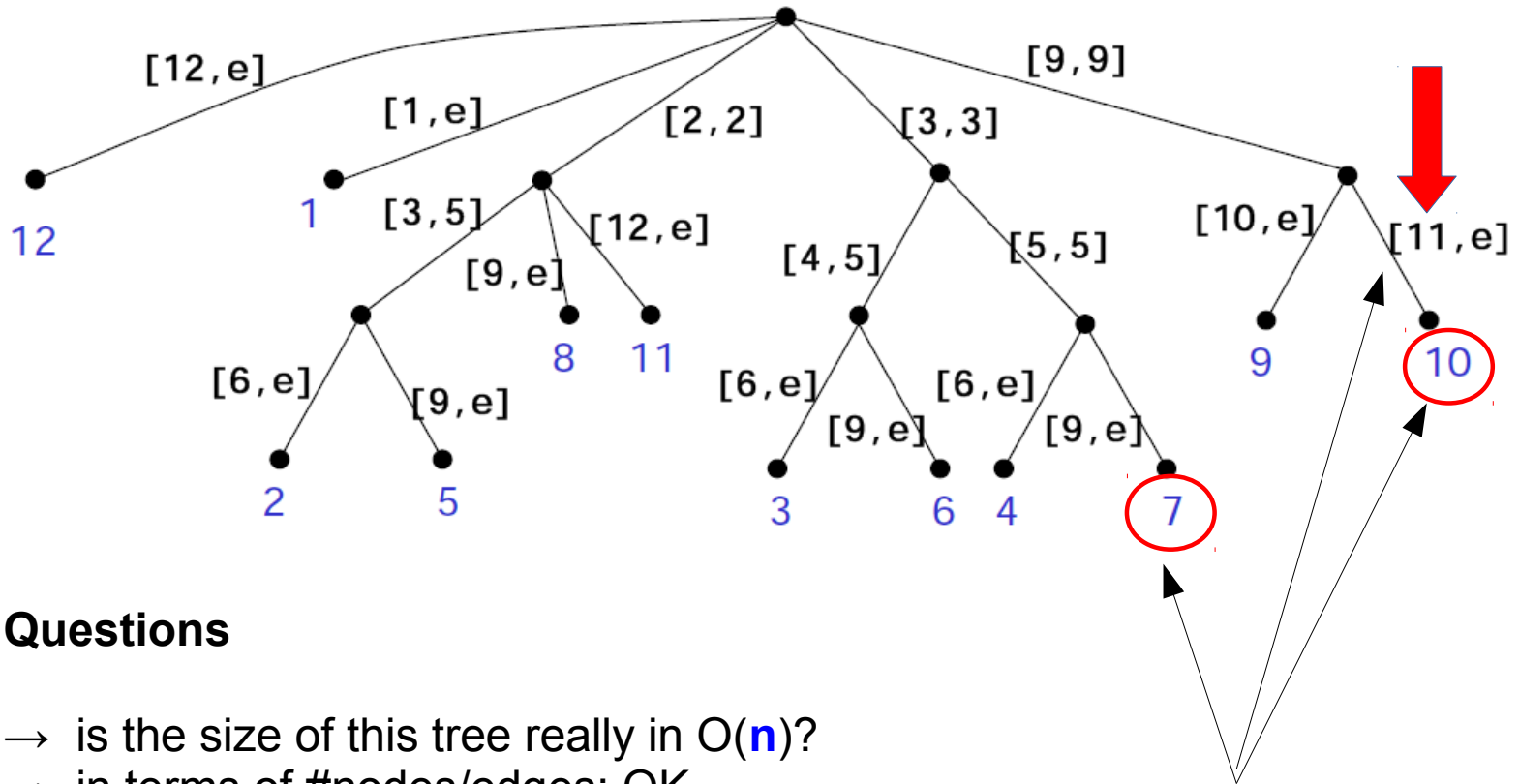
T = mississippi\$



Questions

- is the size of this tree really in $O(n)$?
- in terms of #nodes/edges: OK
- how about the **sizes of labels ??**

T = mississippi\$

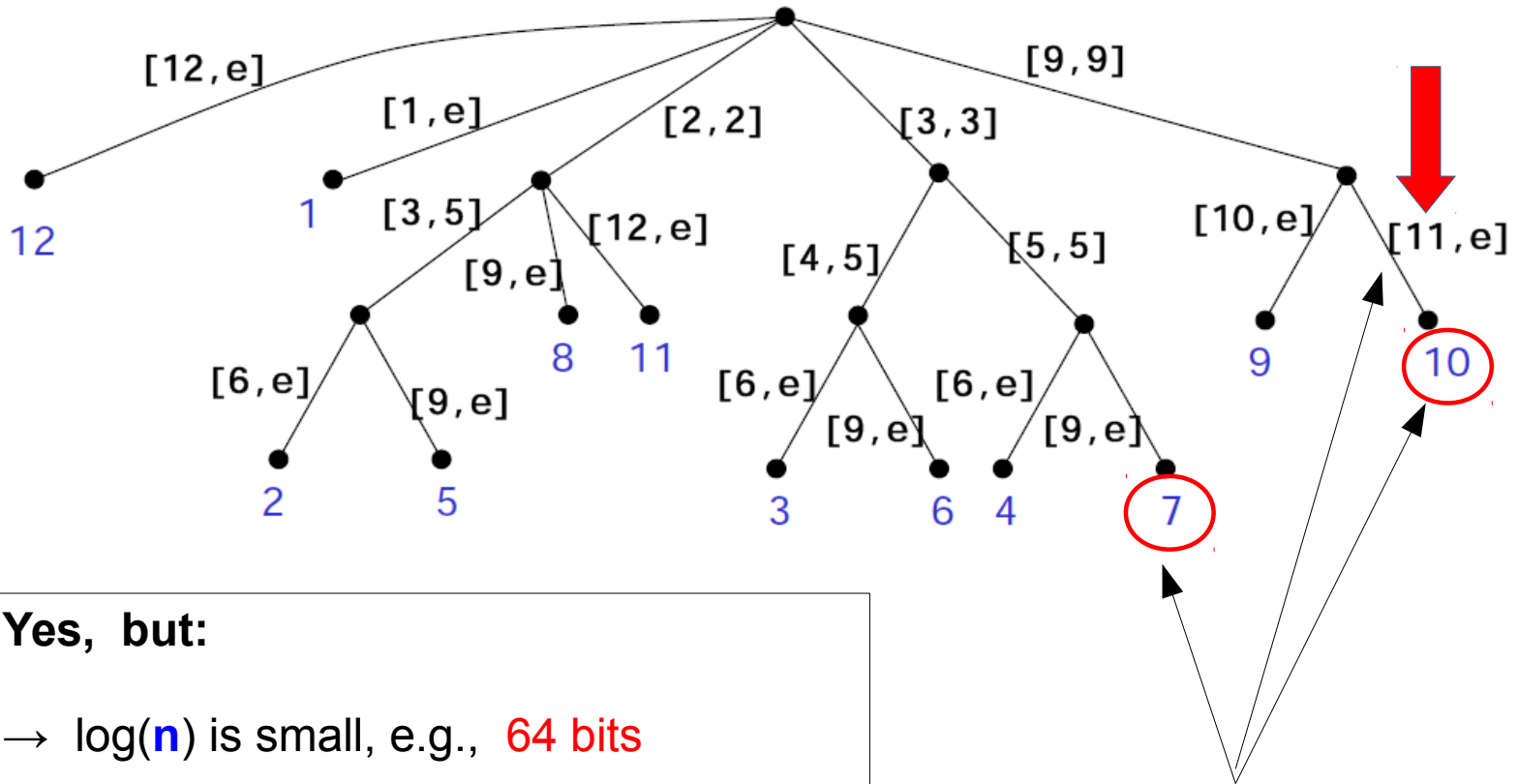


Questions

- is the size of this tree really in $O(n)$?
- in terms of #nodes/edges: OK
- how about the sizes of labels ??

each requires $\log(n)$ bits!?

T = mississippi\$

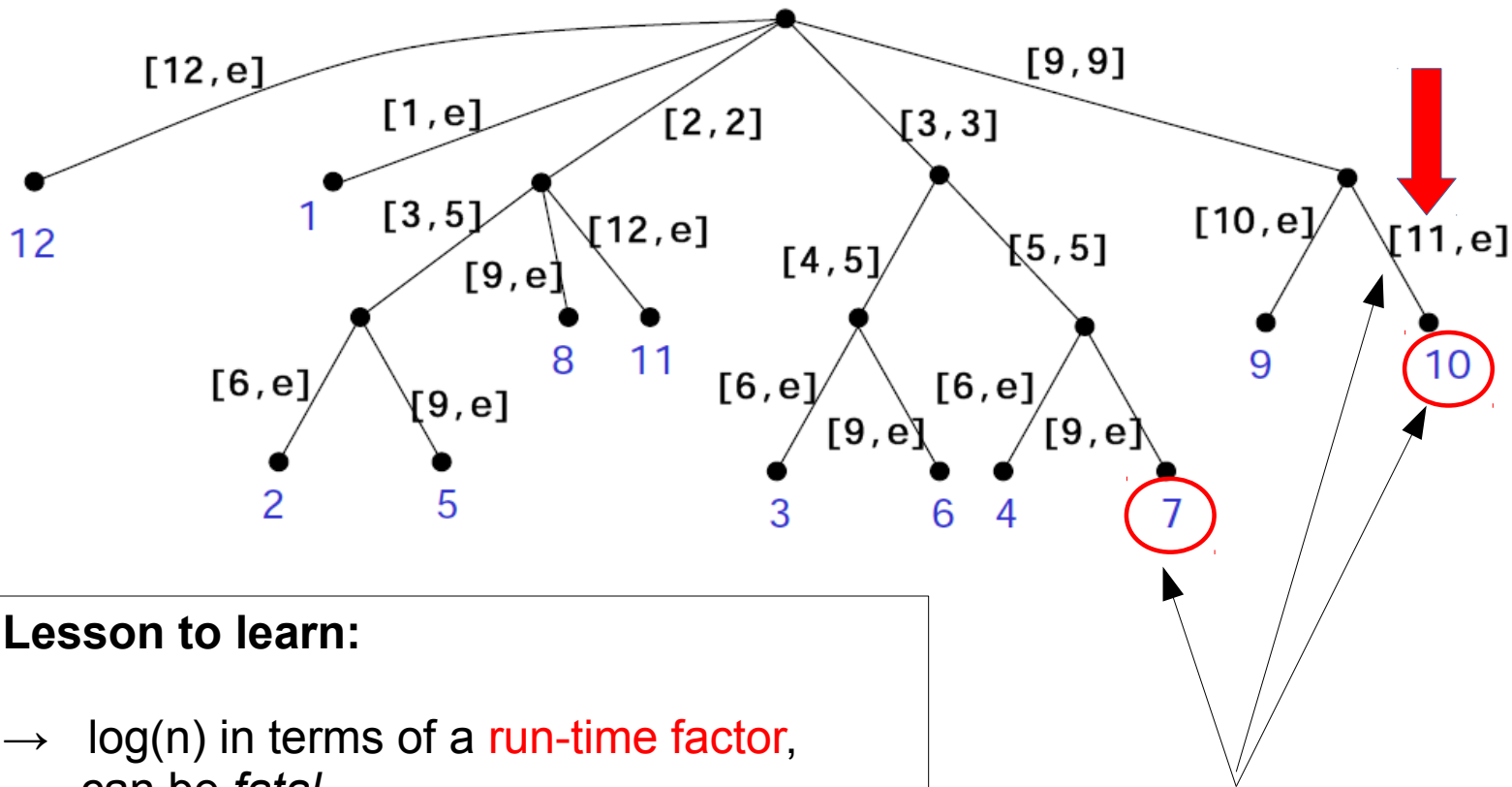


Yes, but:

- $\log(n)$ is small, e.g., 64 bits
- can be considered constant!

each requires $\log(n)$ bits!?

T = mississippi\$

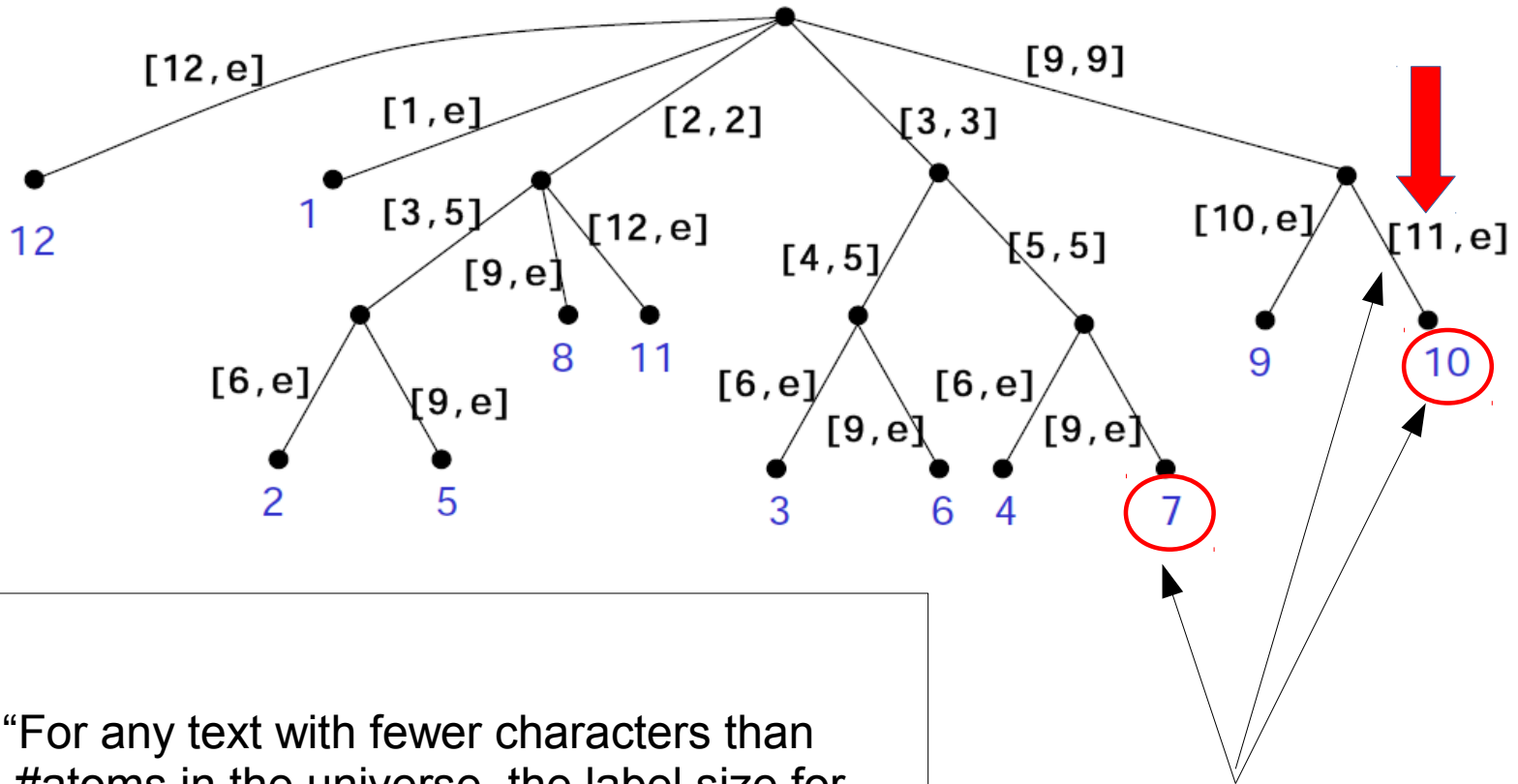


Lesson to learn:

- $\log(n)$ in terms of a **run-time factor**, can be *fatal*
- in terms of a **space-factor**, it is *fine!*
how long will a text be?
 $2^{64}???$

4 / 8 Bytes each is enough!

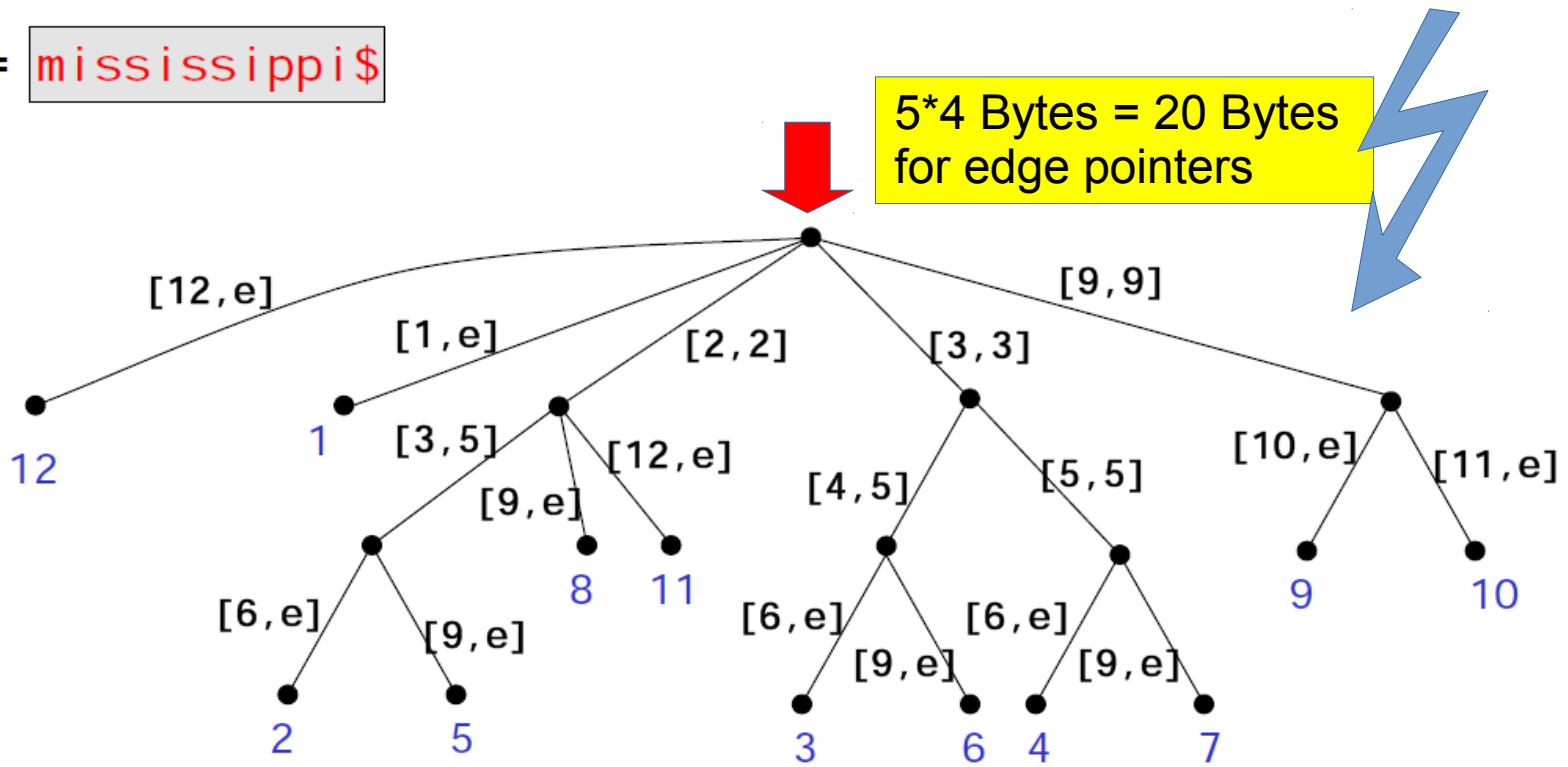
T = mississippi\$



“For any text with fewer characters than #atoms in the universe, the label size for the suffix tree is a constant of **x bits**.. “

x Bits each is enough!

T = mississippi\$



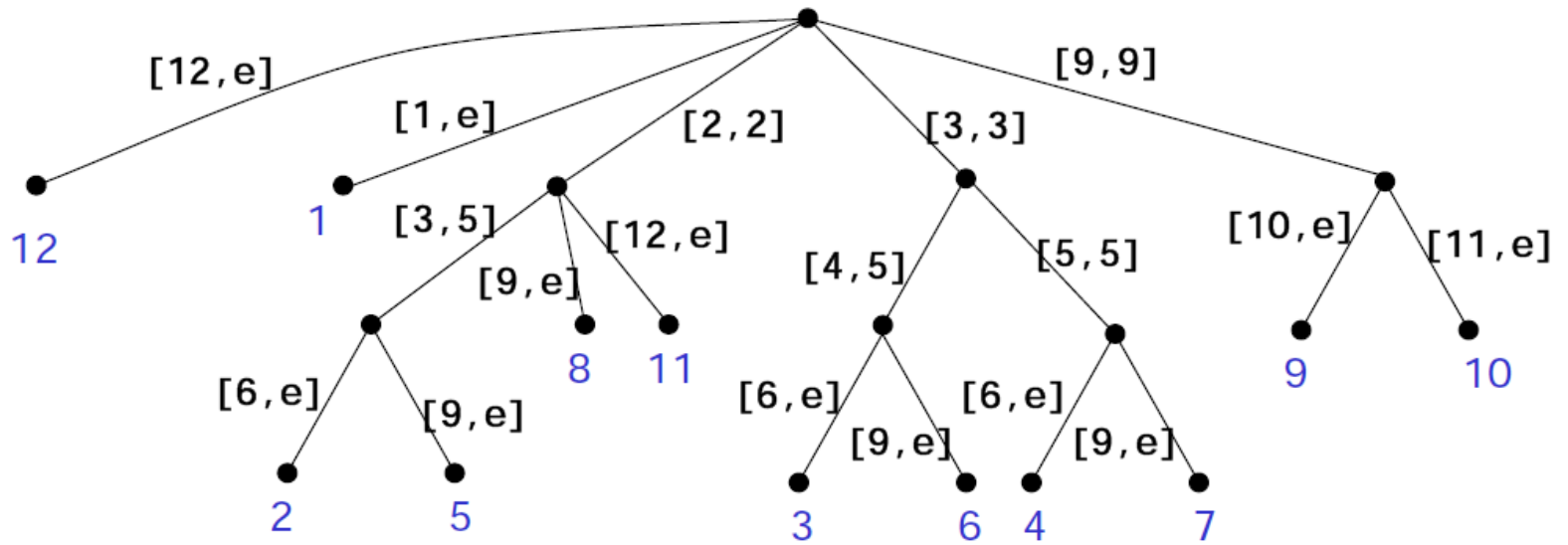
- label size is not an issue
- but, size of **edge-pointers**?
- imagine each edge requires a **32-bit pointer**!!

Actual Space of Suffix Trees

Space for edge-pointers is **problematic**:

→ actual space of suffix tree, ca. **$20|T|$**

→ on commodity hardware, texts of more than 1GB are not doable



→ how to avoid the **huge space needed for edges**?

1. Suffix Array

Definition

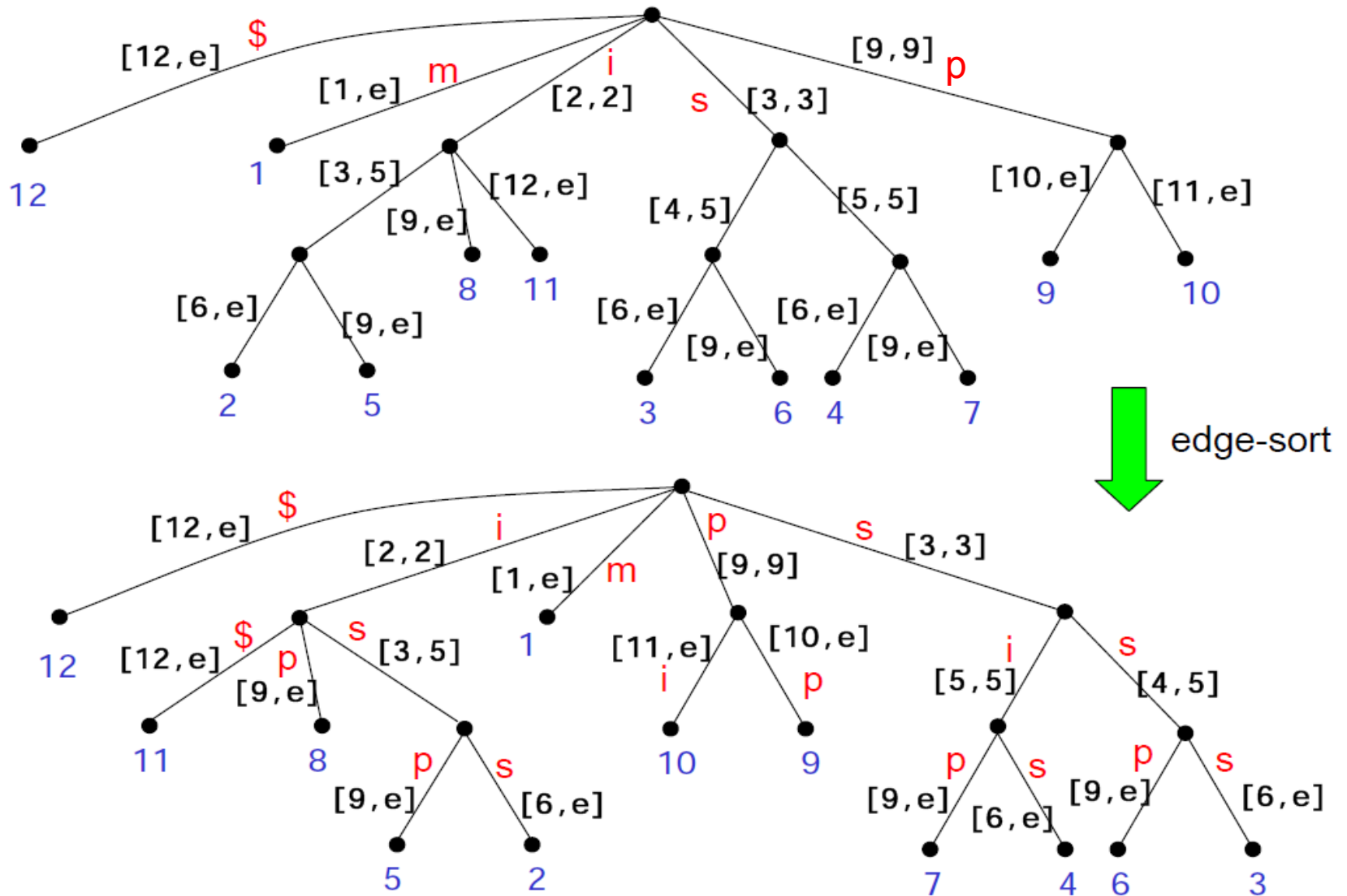
Given text T of length n . For $i=1\dots n$, $SA[k]=i$ if suffix $T[i\dots n]$ is at position k in the lexicographic order T 's suffixes.

1234567890
 $T = \text{mississippi\$}$ Order $\$ < i < m < p < s$

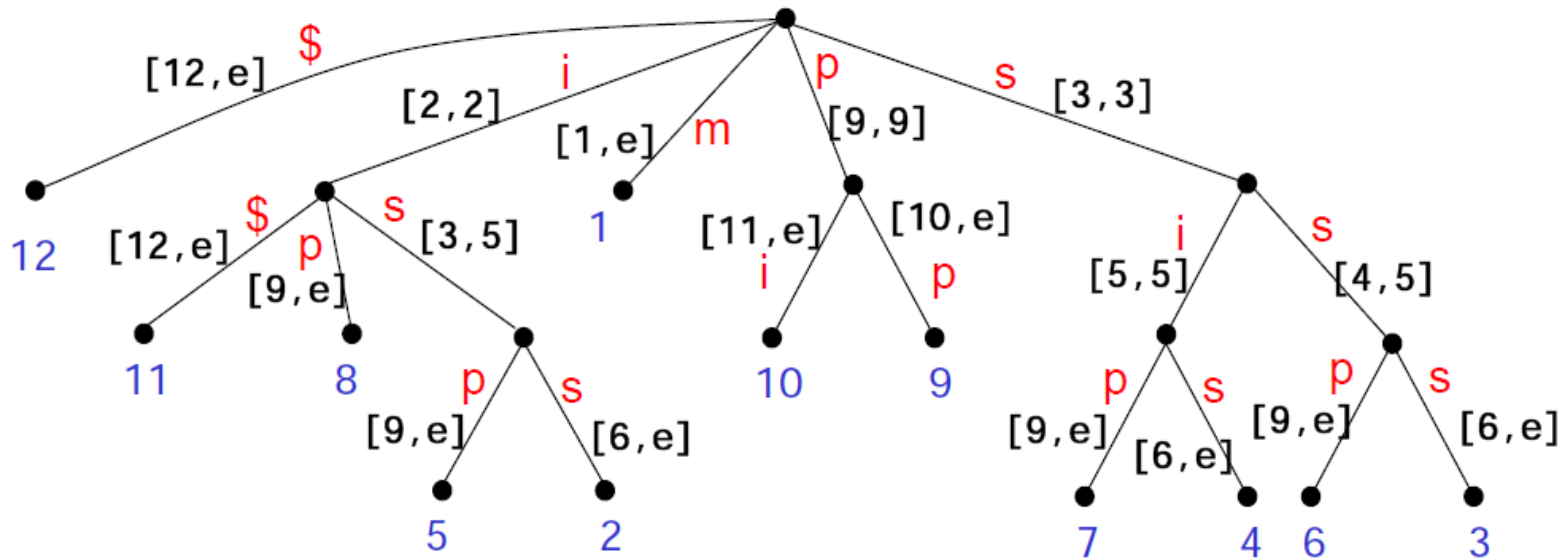
12 $\$$
 11 $i\$$
 8 $ippi\$$
 5 $issippi\$$
 2 $ississippi\$$
 1 $mississippi\$$
 10 $pi\$$
 9 $ppi\$$
 7 $sippi\$$
 4 $sissippi\$$
 6 $ssippi\$$
 3 $ssissippi\$$

$SA(T) = [12, 11, 8, 5, 2, 1, 10, 9, 7, 4, 6, 3]$

Suffix Array Construction



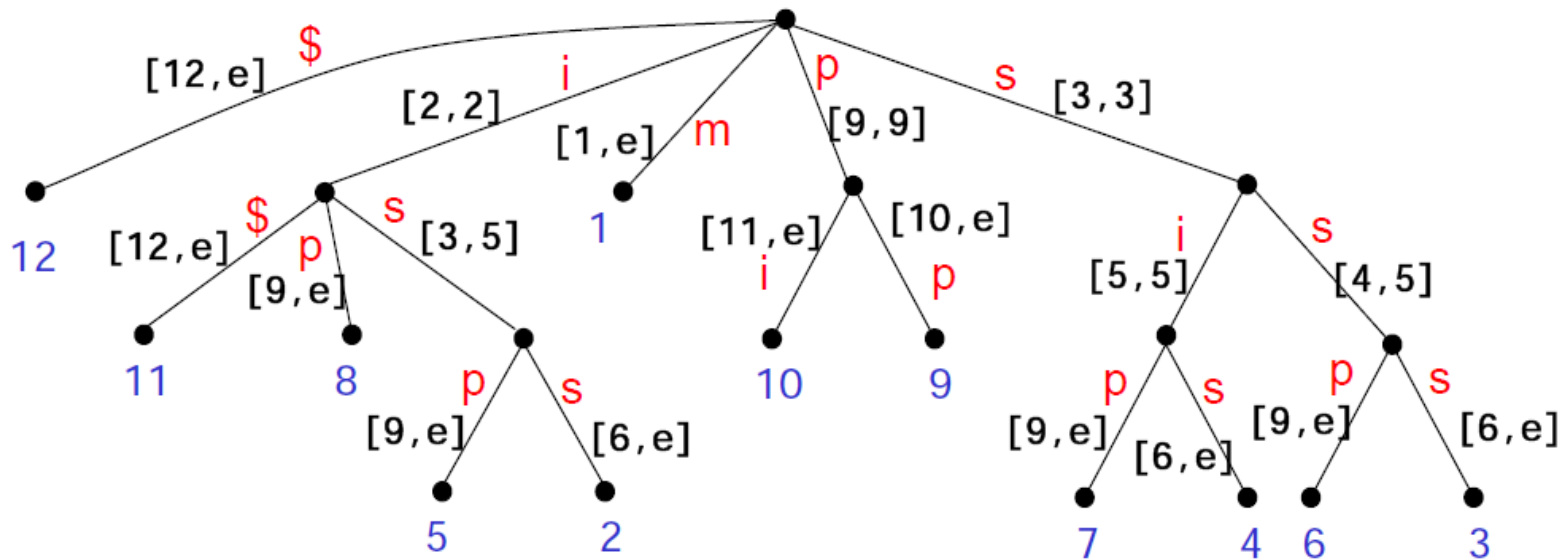
1. Suffix Array



→ read leaves from left-to-right!

SA(T) = [12, 11, 8, 5, 2, 10, 9, 7, 4, 6, 3]

1. Suffix Array



→ read leaves from left-to-right!

SA(T) = [12, 11, 8, 5, 2, 10, 9, 7, 4, 6, 3]

Theorem

The **suffix array** of T can be constructed in time $O(|T|)$.

Search

Theorem

Using binary search on $SA(T)$, all occurrences of P in T can be located in $O(|P| * \log|T|)$ time.

1234567890

$T = \text{mississippi\$}$

$SA(T) = [12, 11, 8, 5, 2, 1, 10, 9, 7, 4, 6, 3]$

Search for $P = \text{issi}$

all occurren's
consecutive in $SA!$

12	\$
11	i\$
8	ippi\$
5	issippi\$
2	issippi\$
1	mississippi\$
10	pi\$
9	ppi\$
7	sippi\$
4	sissippi\$
6	ssippi\$
3	ssissippi\$

Binary search for start-index:

$L=1, R=|T|=n$

Repeat

$M = \lceil (L+R-1)/2 \rceil$

If $P \leq_{\text{lex}} T[M \dots M+|P|]$ then $R:=M$ else $L:=M$

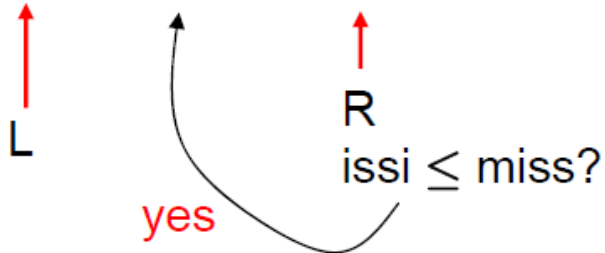
Until M does not change.

Search

1234567890

T = mississippi\$

SA(T) = [12, 11, 8, 5, 2, 1, 10, 9, 7, 4, 6, 3]



12 \$
 11 i\$
 8 ippi\$
 5 issippi\$
 2 ississippi\$
 1 mississippi\$
 10 pi\$
 9 ppi\$
 7 sippi\$
 4 sissippi\$
 6 ssippi\$
 3 ssissippi\$

Binary search for start-index:

$L=1$, $R=|T|=n$

Repeat

$M = \lceil (L+R-1)/2 \rceil$

If $P \leq_{\text{lex}} T[M \dots M+|P|]$ then $R:=M$ else $L:=M$

Until M does not change.

Search

1234567890

T = mississippi\$

SA(T) = [12, 11, 8, 5, 2, 1, 10, 9, 7, 4, 6, 3]

↑ ↑ ↑
L M R

issi ≤ issi

12 \$
11 i\$
8 ippi\$
5 **iss**ippi\$
2 **iss**issippi\$
1 mississippi\$
10 pi\$
9 ppi\$
7 sippi\$
4 sissippi\$
6 ssippi\$
3 ssissippi\$

Binary search for start-index:

L=1, R=|T|=n

Repeat

M = $\lceil (L+R-1)/2 \rceil$

If $P \leq_{\text{lex}} T[M \dots M+|P|]$ then R:=M else L:=M

Until M does not change.

Search

1234567890
 T = mississippi\$

SA(T) = [12, 11, 8, 5, 2, 1, 10, 9, 7, 4, 6, 3]

↑ ↑
 L R

issi ≤ issi

12 \$
 11 i\$
 8 ippi\$
 5 **iss**ippi\$
 2 **iss**iissippi\$
 1 mississippi\$
 10 pi\$
 9 ppi\$
 7 sippi\$
 4 sissippi\$
 6 ssippi\$
 3 ssissippi\$

Binary search for start-index:

L=1, R=|T|=n

Repeat

M = $\lceil (L+R-1)/2 \rceil$

If $P \leq_{\text{lex}} T[M \dots M+|P|]$ then R:=M else L:=M

Until M does not change.

Search

1234567890

T = mississippi\$

SA(T) = [12, 11, 8, 5, 2, 1, 10, 9, 7, 4, 6, 3]

R

issi ≤ issi

M does not change.

→ Return R

Binary search for start-index:

L=1, R=|T|=n

Repeat

M = $\lceil (L+R-1)/2 \rceil$

If $P \leq_{\text{lex}} T[M \dots M+|P|]$ then R:=M else L:=M

Until M does not change.

12 \$
 11 i\$
 8 ippi\$
 5 issippi\$
 2 ississippi\$
 1 mississippi\$
 10 pi\$
 9 ppi\$
 7 sippi\$
 4 sissippi\$
 6 ssippi\$
 3 ssissippi\$

Search

1234567890
 T = mississippi\$

Start-index

SA(T) = [12, 11, 8, 5, 2, 1, 10, 9, 7, 4, 6, 3]

12 \$
 11 i\$
 8 ippi\$
 5 **i**ssippi\$
 2 **i**ssissippi\$
 1 mississippi\$
 10 pi\$
 9 ppi\$
 7 sippi\$
 4 sissippi\$
 6 ssippi\$
 3 sissippi\$

Binary search for **end**-index:

$L=1$, $R=|T|=n$

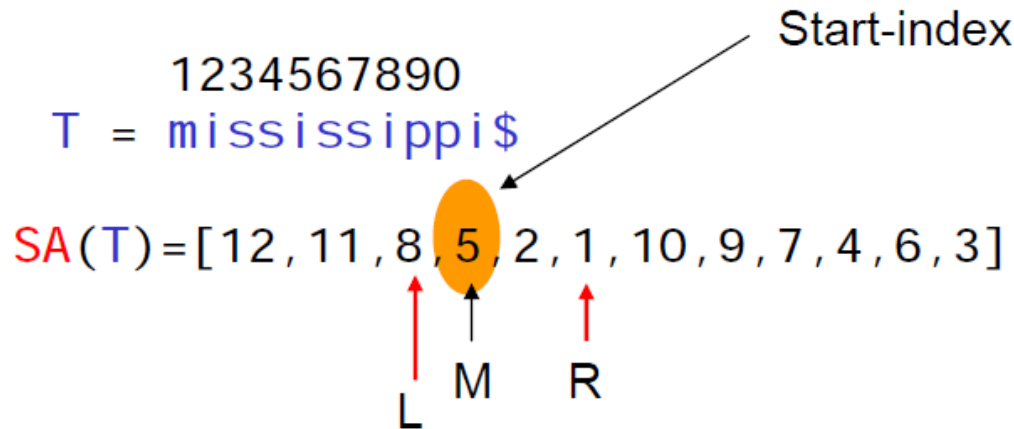
Repeat

$M = \lfloor (L+R-1)/2 \rfloor$

If $P <_{\text{lex}} T[M \dots M+|P|]$ then $R:=M$ else $L:=M$

Until M does not change.

Search



not(issi < issi)
 →

12 \$
 11 i\$
 8 ippi\$
 5 **issippi**\$
 2 **issippi**ssippi\$
 1 mississippi\$
 10 pi\$
 9 ppi\$
 7 sippi\$
 4 sissippi\$
 6 sssippi\$
 3 sssissippi\$

Binary search for **end**-index:

L=1, R=|T|=n

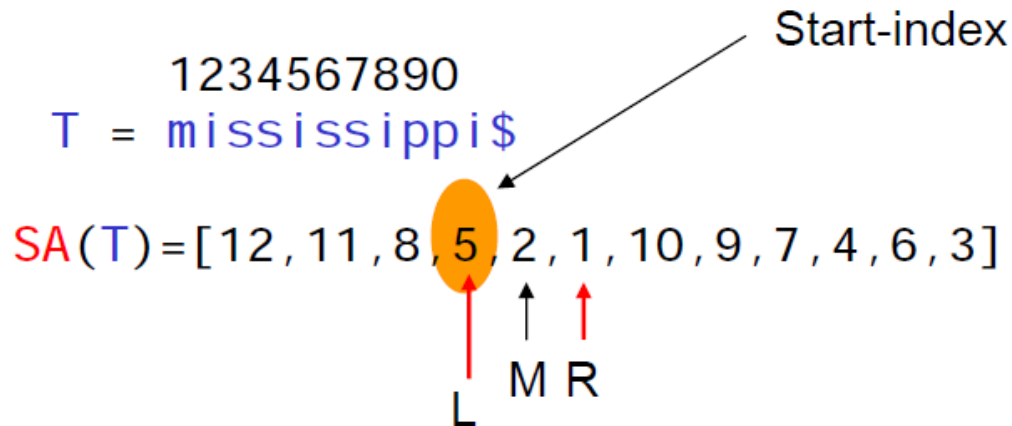
Repeat

M = $\lceil (L+R-1)/2 \rceil$

If $P <_{\text{lex}} T[M \dots M+|P|]$ then R:=M else L:=M

Until M does not change.

Search



not($iss_i < iss_i$)
 $\rightarrow L := M$

12 \$
 11 i\$
 8 ippi\$
 5 issippi\$
 2 ississippi\$
 1 mississippi\$
 10 pi\$
 9 ppi\$
 7 sippi\$
 4 sissippi\$
 6 ssippi\$
 3 ssissippi\$

Binary search for **end**-index:

$L=1, R=|T|=n$

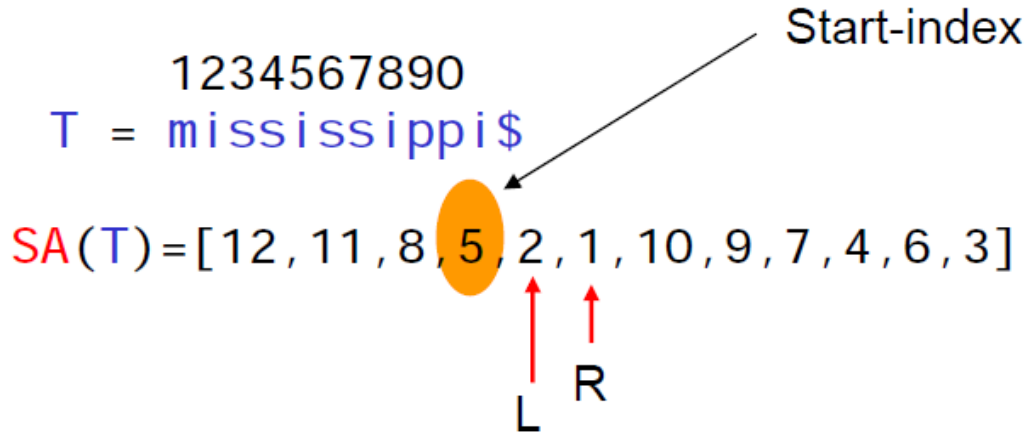
Repeat

$M = \lceil (L+R-1)/2 \rceil$

If $P <_{lex} T[M \dots M+|P|]$ then $R := M$ else $L := M$

Until M does not change.

Search



M does not change.
 → Return L

Binary search for **end**-index:

$L=1$, $R=|T|=n$

Repeat

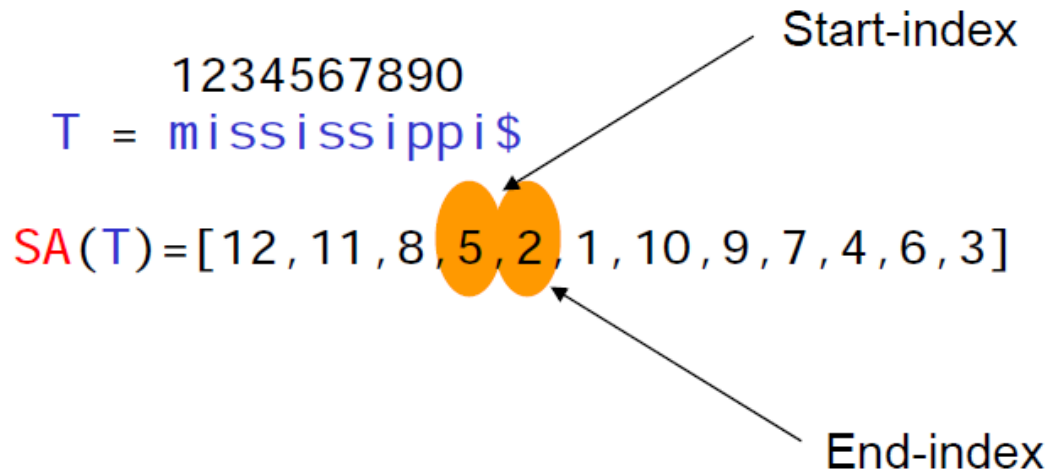
$M = \lceil (L+R-1)/2 \rceil$

If $P <_{lex} T[M\dots M+|P|]$ then $R:=M$ else $L:=M$

Until M does not change.

12 \$
 11 i\$
 8 ippi\$
 5 **iss**ippi\$
 2 **iss**issippi\$
 1 mississippi\$
 10 pi\$
 9 ppi\$
 7 sippi\$
 4 sissippi\$
 6 ssiippi\$
 3 ssiissippi\$

Search



12 \$
 11 i\$
 8 ippi\$
 5 **i**ssippi\$
 2 **i**ssissippi\$
 1 mississippi\$
 10 pi\$
 9 ppi\$
 7 sippi\$
 4 sissippi\$
 6 ssippi\$
 3 ssissippi\$

Binary search for **end**-index:

$L=1, R=|T|=n$

Repeat

$M = \lceil (L+R-1)/2 \rceil$

If $P <_{lex} T[M...M+|P|]$ then $R:=M$ else $L:=M$

Until M does not change.

Search

Theorem

Using binary search on $SA(T)$, all occurrences of P in T can be located in $O(|P| * \log|T|)$ time.

Note

This is a pessimistic bound!

We *almost never* need $O(|P|)$ time for one lexicographic comparison!

On random strings, this should run in $O(|P| + \log|T|)$ time.



Search

Theorem

Using binary search on $SA(T)$, all occurrences of P in T can be located in $O(|P| * \log|T|)$ time.

Note

This is a pessimistic bound!

We *almost never* need $O(|P|)$ time for one lexicographic comparison!

On random strings, this should run in $O(|P| + \log|T|)$ time.

→ $O(|P| + \log|T|)$ in practise, using a simple trick

→ $O(|P| + \log|T|)$ guaranteed, using **LCP-array**

LCP(k,j) = longest common prefix of $T[SA[k]...]$
and $T[SA[j]...]$

1. Suffix Arrays

- much more space efficient than Suffix Tree
 - used in practise (suffix tree more used in theory)
-

→ Suffix Array Construction, without Suffix Trees?

[[Linear Work Suffix Array Construction](#),
J. Kärkkäinen, Sanders, Burkhardt,
Journal of the ACM, 2006]

→ See also:

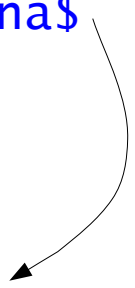
[[A taxonomy of suffix array construction algorithms](#),
S. J. Puglisi, W. F. Smyth, A. Turpin,
ACM Computing Surveys 39, 2007]

2. Burrows-Wheeler Transform

T = banana\$

banana\$
\$banana
a\$banan
na\$bana
ana\$ban
nana\$ba
anana\$b

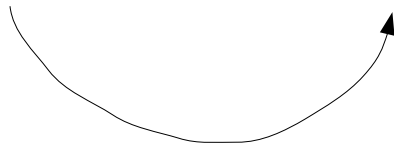
generate all
cyclic shifts of T



2. Burrows-Wheeler Transform

T = banana\$

banana\$	\$banana
\$banana	a\$banan
a\$banan	ana\$ban
na\$bana	anana\$b
ana\$ban	banana\$
nana\$ba	na\$bana
anana\$b	nana\$ba



sort them
lexicographically

\$ < a < b < c < . . .

2. Burrows-Wheeler Transform

T = banana\$

banana\$	\$banana
\$banana	a\$banan
a\$banan	ana\$ban
na\$bana	anana\$b
ana\$ban	banana\$
nana\$ba	na\$bana
anana\$b	nana\$ba

What information is captured by the **first column**?

sort them
lexicographically

\$ < a < b < c < . . .

2. Burrows-Wheeler Transform

T = banana\$

banana\$	\$banana
\$banana	a\$banan
a\$banan	ana\$ban
na\$bana	anana\$b
ana\$ban	banana\$
nana\$ba	na\$bana
anana\$b	nana\$ba

sort them
lexicographically

\$ < a < b < c < . . .

What information is captured by the **first column**?

→ sorted #occ of each letter:

- one time "\$"
- three times "a"
- one time "b"
- two times "n"

2. Burrows-Wheeler Transform

T = banana\$

banana\$	\$banana
\$banana	a\$banan
a\$banan	ana\$ban
na\$bana	anana\$b
ana\$ban	banana\$
nana\$ba	na\$bana
anana\$b	nana\$ba

sort them
lexicographically

\$ < a < b < c < . . .

What information is captured by the **first column**?

→ sorted #occ of each letter:

- one time "\$"
- three times "a"
- one time "b"
- two times "n"

Can you retrieve the original text T, given only the first column?

2. Burrows-Wheeler Transform

T = banana\$

banana\$	\$banana
\$banana	a\$banan
a\$banan	ana\$ban
na\$bana	anana\$b
ana\$ban	banana\$
nana\$ba	na\$bana
anana\$b	nana\$ba

sort them
lexicographically

\$ < a < b < c < . . .

What information is captured by the **first column**?

→ sorted #occ of each letter:

- one time “\$”
- three times “a”
- one time “b”
- two times “n”

Can you retrieve the original text T, given only the first column?

Of course not!!

2. Burrows-Wheeler Transform

$T = \text{banana}\$$

banana\$	\$banana
\$banana	a\$banan
a\$banan	ana\$ban
na\$bana	anana\$b
ana\$ban	banana\$
nana\$ba	na\$bana
anana\$b	nana\$ba

sort them
lexicographically

\$ < a < b < c < . . .

What information is captured by the **first column**?

→ sorted #occ of each letter:

- one time "\$"
- three times "a"
- one time "b"
- two times "n"

Can you retrieve the original text T , given only the first column?

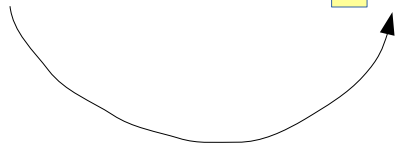
Note

→ each column contains the **same letters** (\$, 3*a, b, 2*n)

2. Burrows-Wheeler Transform

T = banana\$

banana\$	\$banana
\$banana	a\$banan
a\$banan	ana\$ban
na\$bana	anana\$b
ana\$ban	banana\$
nana\$ba	na\$bana
anana\$b	nana\$ba



sort them
lexicographically

What information is captured
by the **second column**?

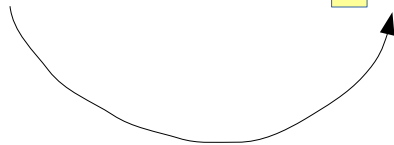
→ “sorting with respect to considering
one previous letter”

\$ < a < b < c < . . .

2. Burrows-Wheeler Transform

$T = \text{banana}\$$

banana\$	\$banana
\$banana	a\$banan
a\$banan	ana\$ban
na\$bana	anana\$b
ana\$ban	banana\$
nana\$ba	na\$bana
anana\$b	nana\$ba



sort them
lexicographically

$\$ < a < b < c < \dots$

What information is captured by the **second column**?

→ “sorting with respect to considering one previous letter”

Can you retrieve the original T from the second column only?

2. Burrows-Wheeler Transform

T = banana\$

banana\$	\$banana
\$banana	a\$banan
a\$banan	ana\$ban
na\$bana	anana\$b
ana\$ban	banana\$
nana\$ba	na\$bana
anana\$b	nana\$ba

sort them
lexicographically

\$ < a < b < c < . . .

What information is captured by the **second column**?

→ “sorting with respect to considering one previous letter”

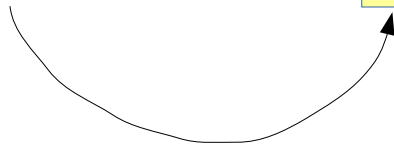
Can you retrieve the original T from the second column only?

Can you retrieve the first column, given the second one?

2. Burrows-Wheeler Transform

T = banana\$

banana\$	\$banana
\$banana	a\$banan
a\$banan	ana\$ban
na\$bana	anana\$b
ana\$ban	banana\$
nana\$ba	na\$bana
anana\$b	nana\$ba



sort them
lexicographically

What information is captured
by the **third column**?

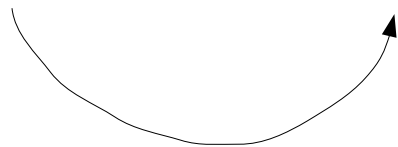
→ “sorting with respect to considering
two previous letters”

\$ < a < b < c < . . .

2. Burrows-Wheeler Transform

T = banana\$

banana\$	\$banana
\$banana	a\$banan
a\$banan	ana\$ban
na\$bana	anana\$b
ana\$ban	banana\$
nana\$ba	na\$bana
anana\$b	nana\$ba



sort them
lexicographically

What information is captured
by the **last column**?

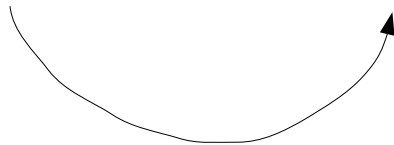
→ “sorting with respect to considering
all previous letters”

\$ < a < b < c < . . .

2. Burrows-Wheeler Transform

T = banana\$

banana\$	\$banana
\$banana	a\$banan
a\$banan	ana\$ban
na\$bana	anana\$b
ana\$ban	banana\$
nana\$ba	na\$bana
anana\$b	nana\$ba



sort them
lexicographically

\$ < a < b < c < . . .

What information is captured
by the **last column**?

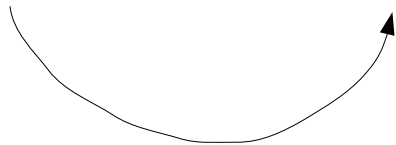
→ “sorting with respect to considering
all previous letters”

Can you retrieve the original text T,
given only the last column?

2. Burrows-Wheeler Transform

T = banana\$

banana\$	\$banana
\$banana	a\$banan
a\$banan	ana\$ban
na\$bana	anana\$b
ana\$ban	banana\$
nana\$ba	na\$bana
anana\$b	nana\$ba



sort them
lexicographically

\$ < a < b < c < . . .

What information is captured
by the **last column**?

→ “sorting with respect to considering
all previous letters”

Can you retrieve the original text T,
given only the last column?

YES, you can!!

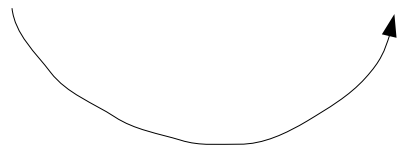
2. Burrows-Wheeler Transform

$T = \text{banana}\$$

banana\$
\$banana
a\$banan
na\$bana
ana\$ban
nana\$ba
anana\$b

\$banana
a\$banan
ana\$ban
anana\$b
banana\$
na\$bana
nana\$ba

Burrows-Wheeler Transform L
of text T



sort them
lexicographically

$\$ < a < b < c < \dots$

2. Burrows-Wheeler Transform

$T = \text{banana}\$$

banana\$	\$banana
\$banana	a\$banan
a\$banan	ana\$ban
na\$bana	anana\$b
ana\$ban	banana\$
nana\$ba	na\$bana
anana\$b	nana\$ba

Burrows-Wheeler Transform L
of text T

Given **last column** L , how can we
reconstruct the original **text** T ?

sort them
lexicographically

$\$ < a < b < c < \dots$

2. Burrows-Wheeler Transform

Naively:

a
n
n
b
\$
a
a

2. Burrows-Wheeler Transform

Naively:

a
n
n
b
\$
a
a

sort

\$
a
a
a
a
b
n
n

first column!

2. Burrows-Wheeler Transform

Naively:

a
n
n
b
\$
a
a

sort

\$
a
a
a
b
n
n

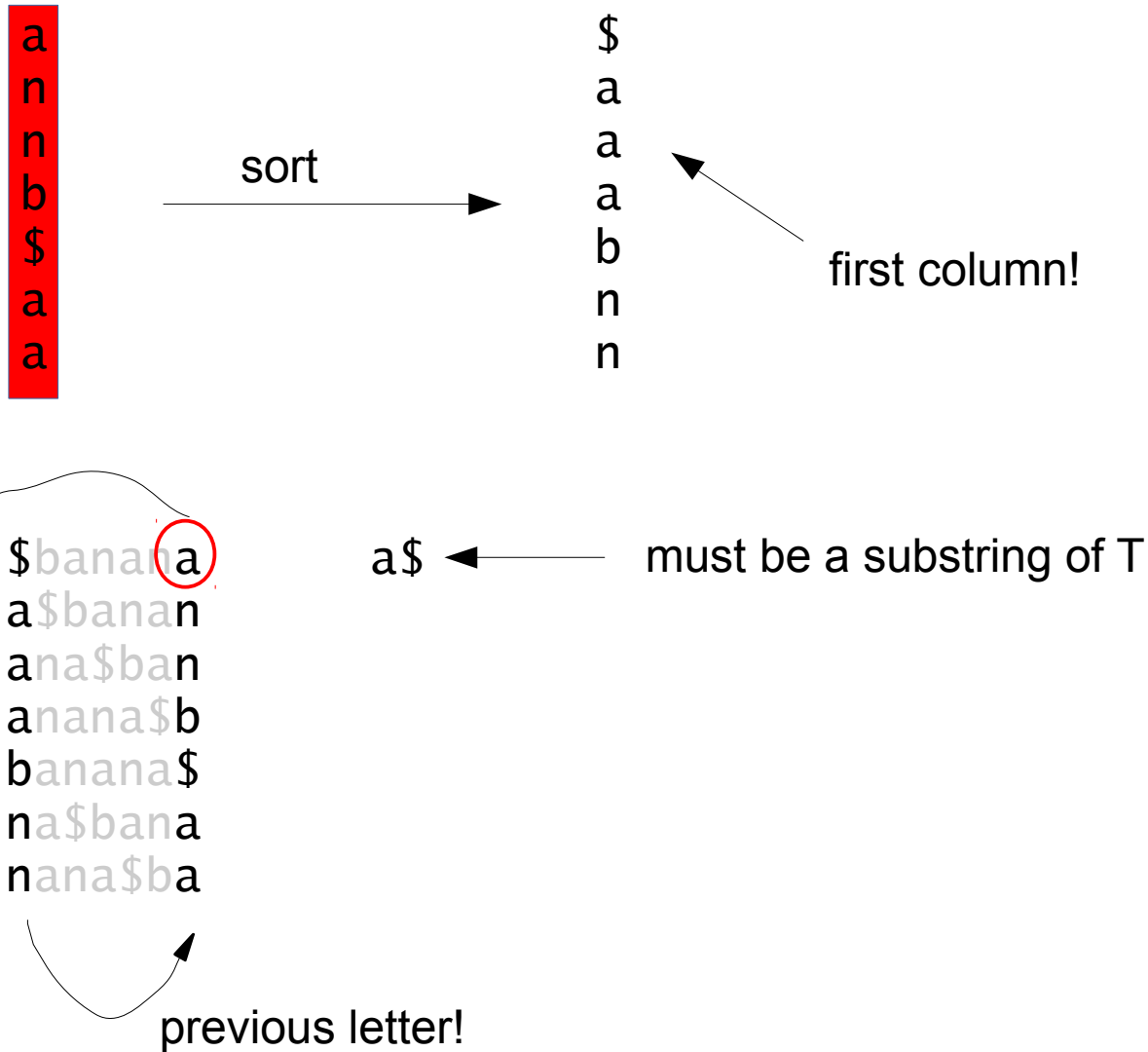
first column!

\$banana
a\$banan
ana\$ban
anana\$b
banana\$
na\$bana
nana\$ba

previous letter!

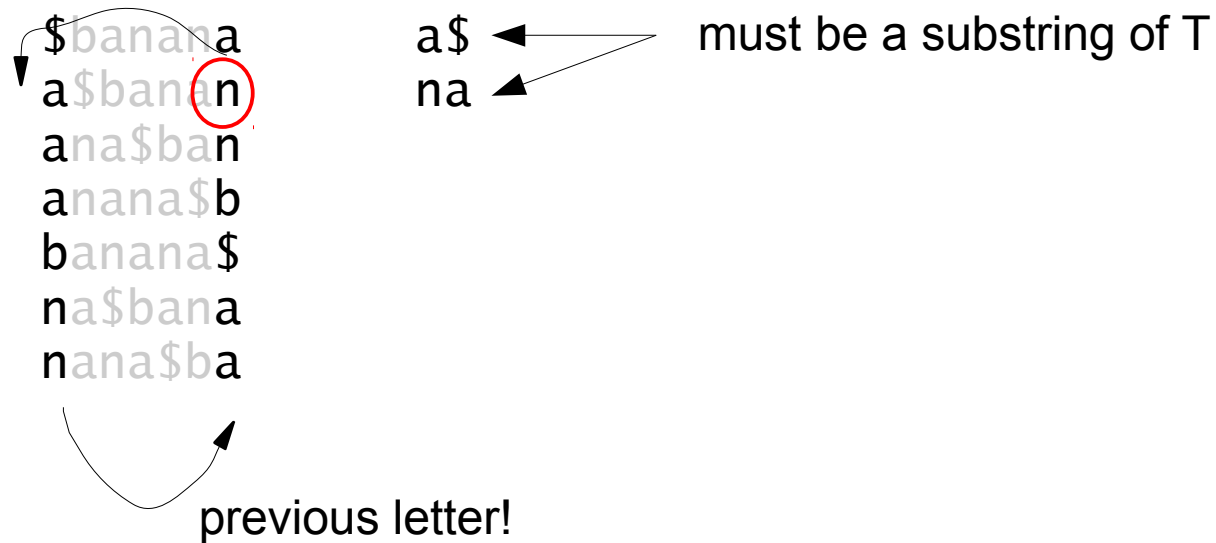
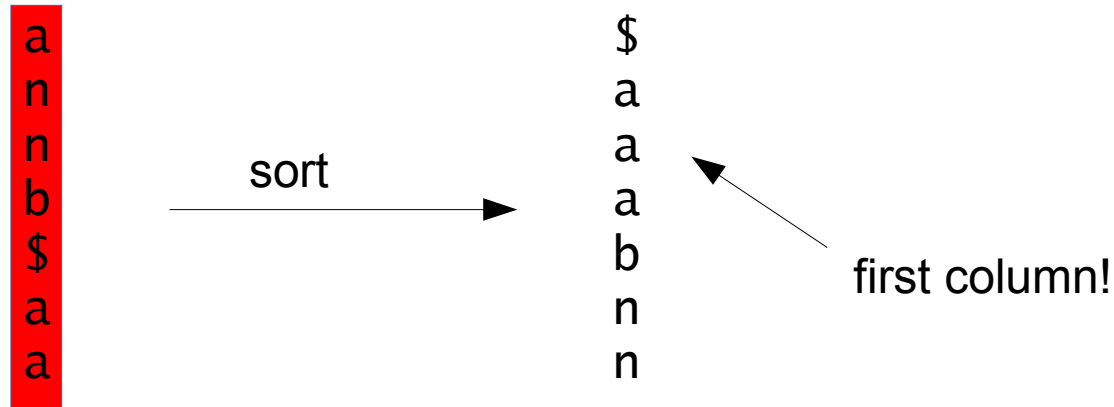
2. Burrows-Wheeler Transform

Naively:



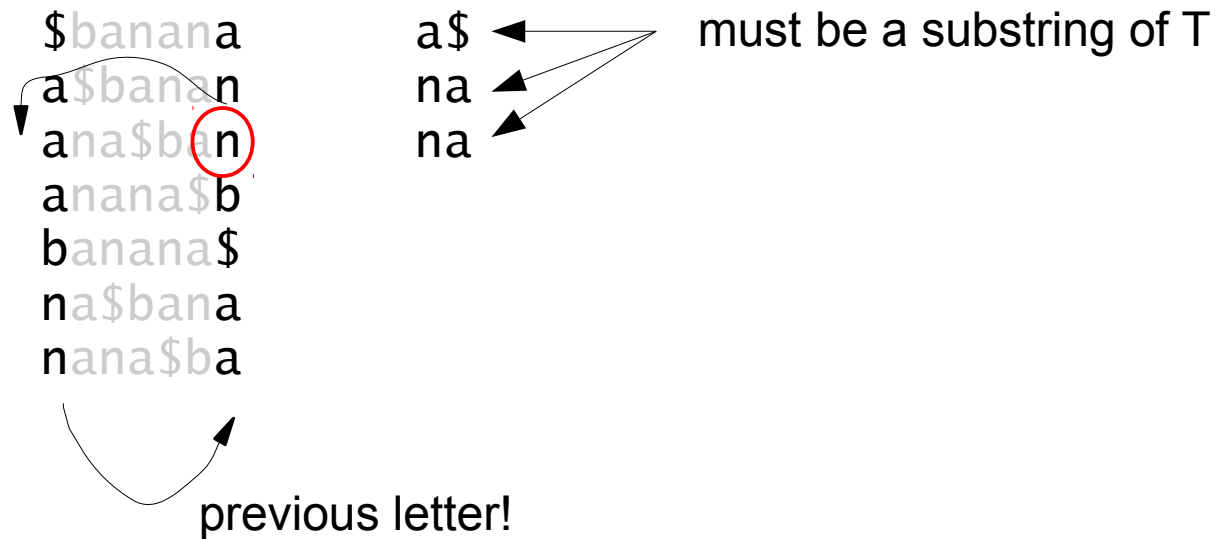
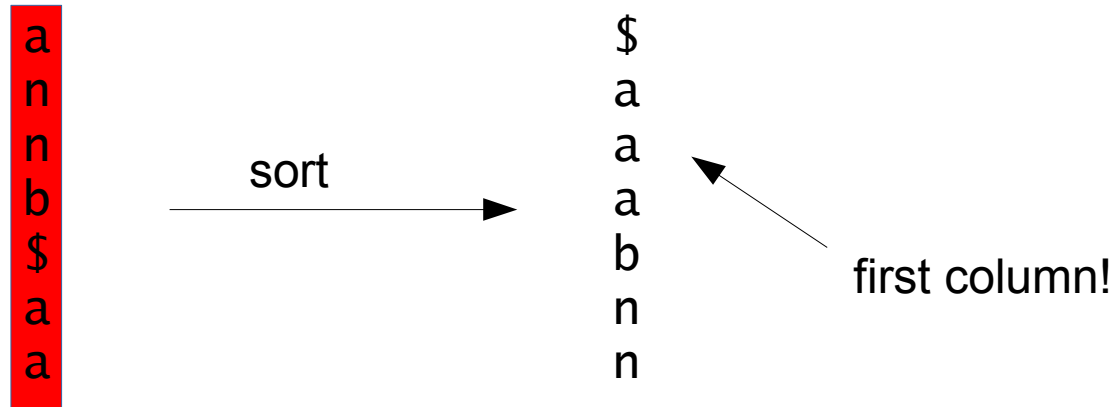
2. Burrows-Wheeler Transform

Naively:



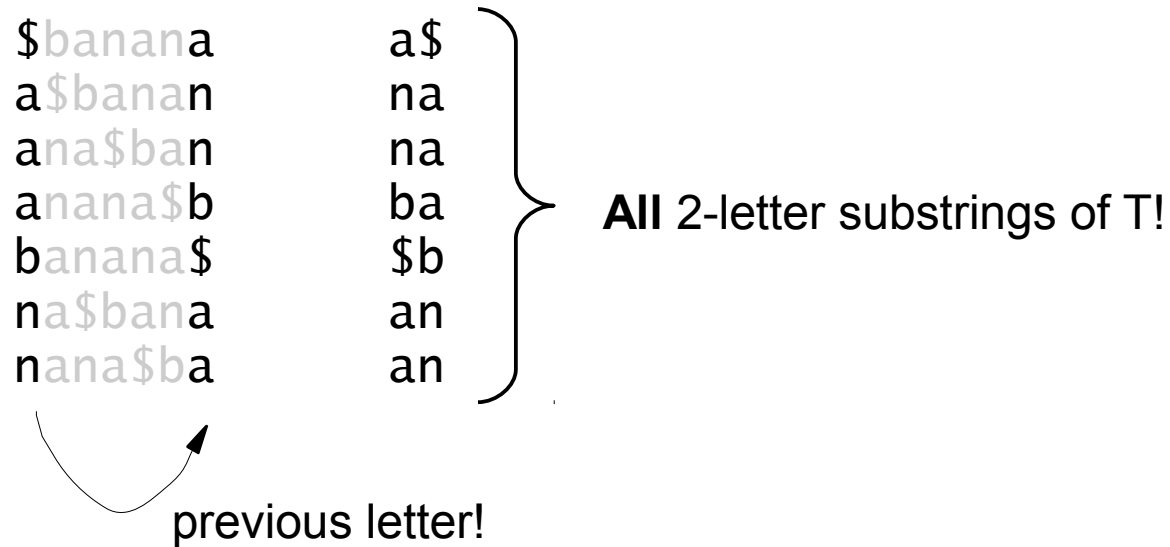
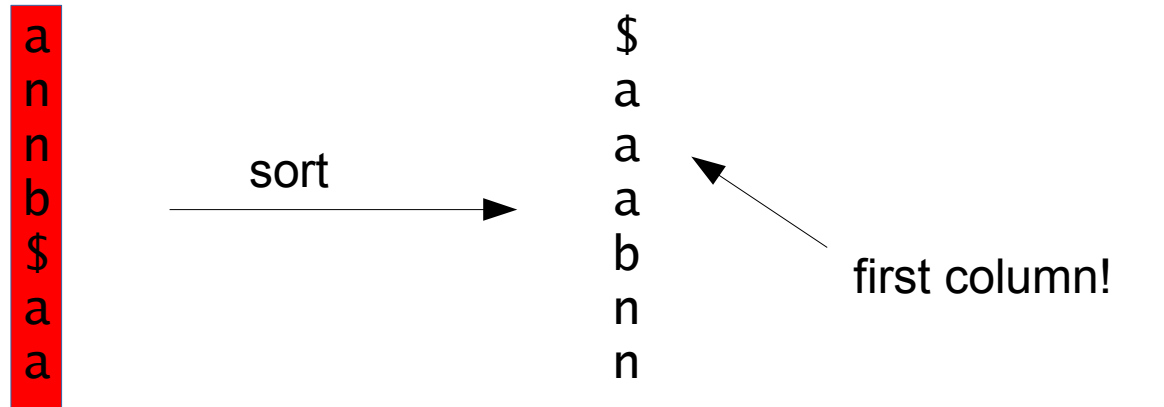
2. Burrows-Wheeler Transform

Naively:



2. Burrows-Wheeler Transform

Naively:



2. Burrows-Wheeler Transform

Naively:

a
n
n
b
\$
a
a

sort

\$
a
a
a
b
n
n

first column!

a\$
na
na
ba
\$b
an
an

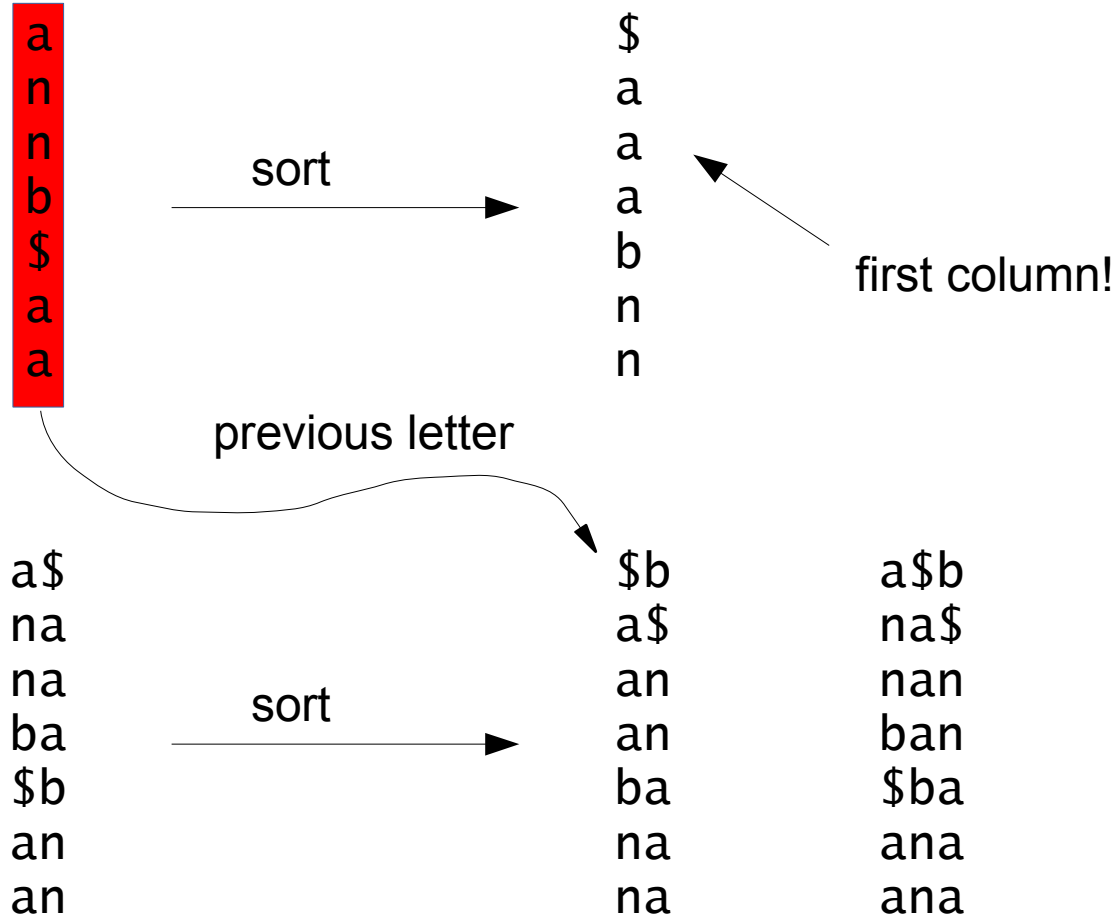
sort

\$b
a\$
an
an
ba
na
na

columns 1 and 2.

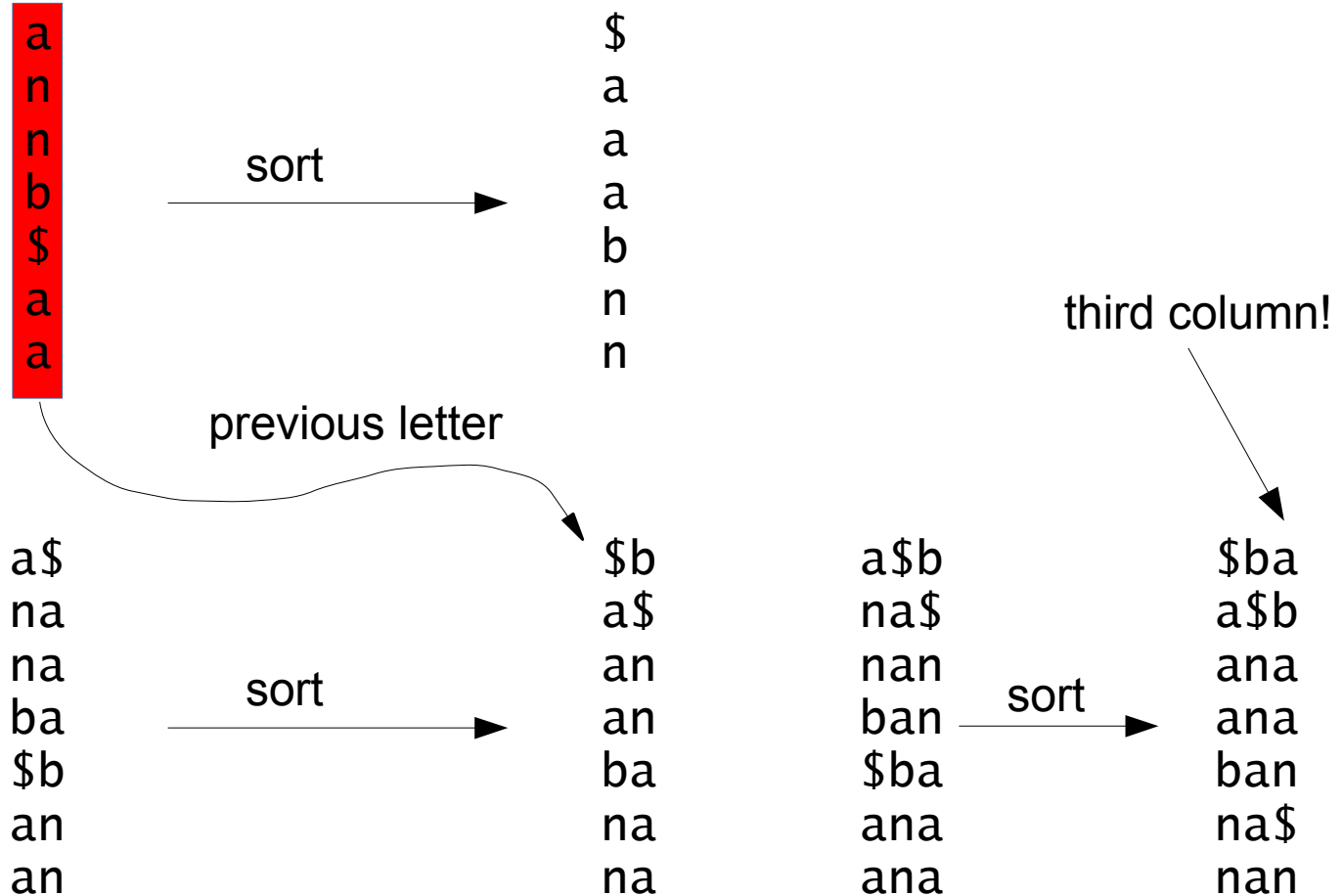
2. Burrows-Wheeler Transform

Naively:



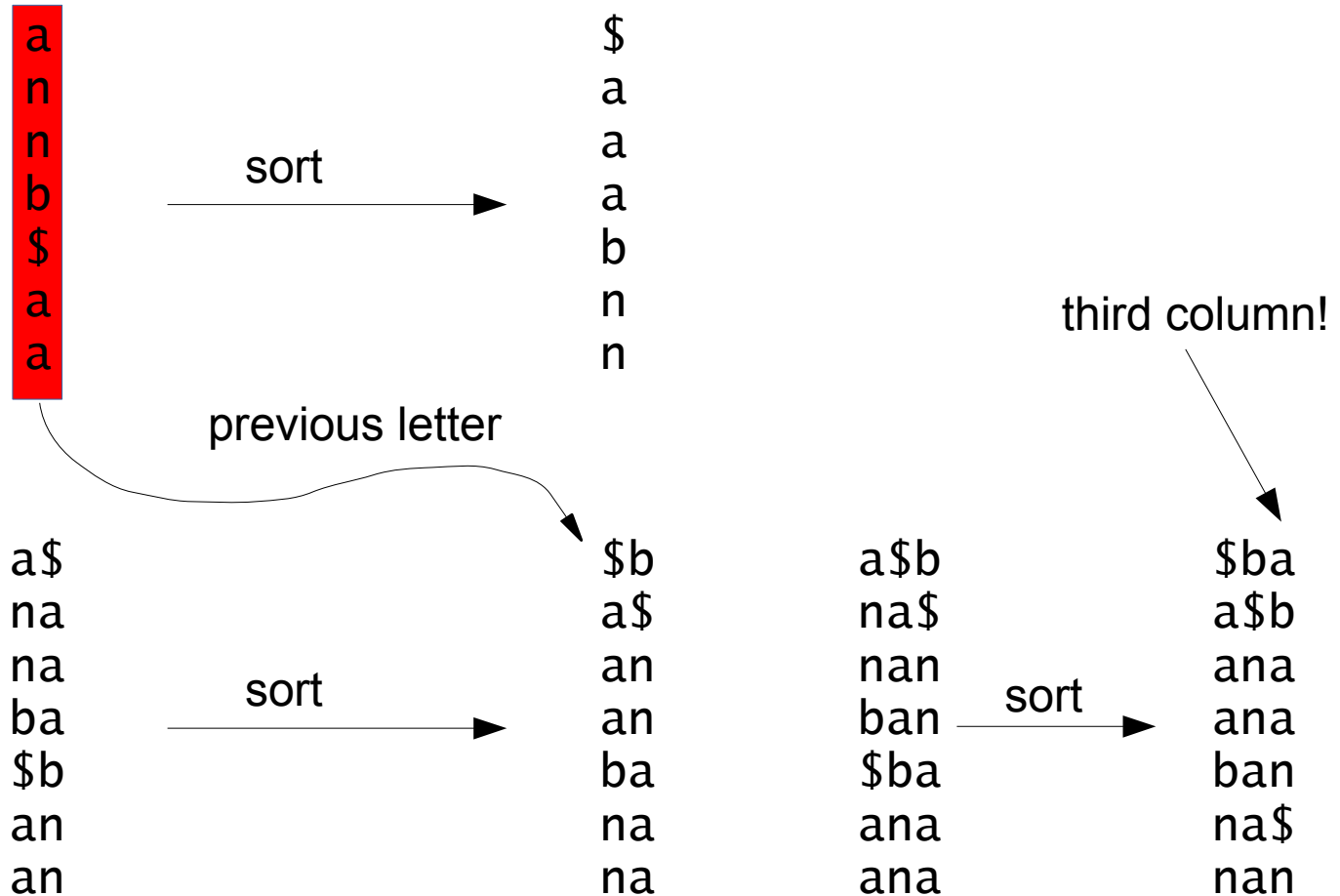
2. Burrows-Wheeler Transform

Naively:



2. Burrows-Wheeler Transform

Naively:



Et cetera

2. Burrows-Wheeler Transform

Naive method: very expensive! (many sortings)

$\text{rank}_b(L, k) = \# \text{occ of } b \text{ in } L, \text{ up to position } k \text{ (excluding } k)$

	1	2	3	4	5	6	7
$L =$	a	n	n	b	\$	a	a

$\text{rank}_n(L, 4) = 2$

2. Burrows-Wheeler Transform

Naive method: very expensive! (many sortings)

$\text{rank}_b(L, k) = \# \text{occ of } b \text{ in } L, \text{ up to position } k \text{ (excluding } k)$

	1	2	3	4	5	6	7
$L =$	a	n	n	b	\$	a	a

$\text{rank}_n(L, 4) = 2$

$\text{rank}_n(L, 3) = 1$

2. Burrows-Wheeler Transform

Naive method: very expensive! (many sortings)

$\text{rank}_b(L, k) = \# \text{occ of } b \text{ in } L, \text{ up to position } k \text{ (excluding } k)$

$$\begin{array}{ccccccc}
 & 1 & 2 & 3 & 4 & 5 & 6 & 7 \\
 L = & a & n & n & b & \$ & a & a
 \end{array}$$

$\text{rank}_n(L, 4) = 2$

$\text{rank}_n(L, 3) = 1$

$\text{rank}_a(L, 7) = 2$

$O(\log |S|)$ time
(after linear time preprocessing of L)

2. Burrows-Wheeler Transform

$\text{rank}_b(L, k) = \# \text{occ of } b \text{ in } L, \text{ up to position } k-1$

$L =$	1	2	3	4	5	6	7	C	1	2	5	6
	a	n	n	b	\$	a	a		\$	a	b	m

1\$
2a
a
a
5b
6n
n

Last-to-Front Mapping

$LF(k) = C[L[k]] + \text{rank}_{L[k]}(L, k)$

first line starting with the letter

\$banana
a\$banan
ana\$ban
anana\$b
banana\$
na\$ban
nana\$ba

▼ second "a", coming from the top
where is the second "a" from top, in the first column?

2. Burrows-Wheeler Transform

$\text{rank}_b(L, k) = \# \text{occ of } b \text{ in } L, \text{ up to position } k-1$

$L =$	1	2	3	4	5	6	7	C	1	2	5	6
	a	n	n	b	\$	a	a		\$	a	b	m

1\$
2a
a
a
5b
6n
n

Last-to-Front Mapping

$\text{LF}(k) = C[L[k]] + \text{rank}_{L[k]}(L, k)$

first line starting with the letter

\$	b	a	n	a	n	a
a	\$	b	a	n	a	n
a	n	a	\$	b	a	n
a	n	a	\$	b	a	n
b	a	n	a	\$	b	a
n	a	\$	b	a	n	a
n	a	\$	b	a	n	a

second "a", coming from the top
where is the second "a" from top, in the first column?

$\text{LF}(6) = C[L[6]] + \text{rank}_a(L, 6) = C["a"] + 1 = 2 + 1 = 3$

Decoding

$\text{rank}_b(L, k) = \# \text{occ of } b \text{ in } L, \text{ up to position } k-1$

$L =$	1	2	3	4	5	6	7	C	1	2	5	6
	a	n	n	b	\$	a	a		\$	a	b	m

Last-to-Front Mapping

$\text{LF}(k) = C[L[k]] + \text{rank}_{L[k]}(L, k)$

first line starting with the letter

previous letter!

\$banana
 a\$anan
 ana\$ban
 anana\$b
 banana\$
 na\$bana
 nana\$ba

→ start with \$ (position 5)

→ $\text{LF}(5) = 1$

→ $L[1] = a$

[current decoding: "a\$"]

Decoding

$\text{rank}_b(L, k) = \# \text{occ of } b \text{ in } L, \text{ up to position } k-1$

$L =$	1	2	3	4	5	6	7	C	1	2	5	6
	a	n	n	b	\$	a	a		\$	a	b	m

first line starting with the letter

Last-to-Front Mapping

$\text{LF}(k) = C[L[k]] + \text{rank}_{L[k]}(L, k)$

\$banana
 a\$banan
 ana\$ban
 anana\$b
 banana\$
 na\$bana
 nana\$ba

→ start with \$ (position 5)

→ $\text{LF}(5) = 1$

→ $L[1] = a$ [current decoding: "a\$"]

→ $\text{LF}(1) = 2$

→ $L[2] = n$ ["ba\$"]

Decoding

$\text{rank}_b(L, k) = \# \text{occ of } b \text{ in } L, \text{ up to position } k-1$

$L =$	1	2	3	4	5	6	7	$\$$	a	b	m	
	a	n	n	b	$\$$	a	a	C	1	2	5	6

first line starting with the letter

Last-to-Front Mapping

$LF(k) = C[L[k]] + \text{rank}_{L[k]}(L, k)$

\$banana
a\$banan
ana\$ban
anana\$b
banana\$
na\$banan
nana\$ba

→ start with $\$$ (position 5)

→ $LF(5) = 1$

→ $L[1] = a$ [current decoding: "a\$"]

→ $LF(1) = 2$

→ $L[2] = n$ ["ba\$ "]

→ $LF(2) = 6$

→ $L(6) = "a"$ ["ana\$ "]

Decoding

$\text{rank}_b(L, k) = \# \text{occ of } b \text{ in } L, \text{ up to position } k-1$

$L =$	1	2	3	4	5	6	7	$\$$	a	b	m	
	a	n	n	b	$\$$	a	a	C	1	2	5	6

first line starting with the letter

Last-to-Front Mapping

$$\text{LF}(k) = C[L[k]] + \text{rank}_{L[k]}(L, k)$$

\$banana
a\$banan
ana\$ban
anana\$b
banana\$
na\$banan
nana\$ba

→ start with \$ (position 5)

→ $\text{LF}(5) = 1$

→ $L[1] = a$ [current decoding: "a\$"]

→ $\text{LF}(1) = 2, L[2] = n$ ["ba\$"]

→ $\text{LF}(2) = 6, L(6) = "a"$ ["ana\$"]

→ $\text{LF}(6) = 3, L(3) = "n"$ ["nana\$"]

Decoding

$\text{rank}_b(L, k) = \# \text{occ of } b \text{ in } L, \text{ up to position } k-1$

$L =$	1	2	3	4	5	6	7	C	1	2	5	6
	a	n	n	b	\$	a	a		\$	a	b	m

first line starting with the letter

Last-to-Front Mapping

$\text{LF}(k) = C[L[k]] + \text{rank}_{L[k]}(L, k)$

\$banana
a\$banan
ana\$ban
anana\$b
banana\$
na\$banan
nana\$ba

→ start with \$ (position 5)

→ $\text{LF}(5) = 1$

→ $L[1] = a$ [current decoding: "a\$"]

→ $\text{LF}(1) = 2, L[2] = n$ ["ba\$"]

→ $\text{LF}(2) = 6, L(6) = "a"$ ["ana\$"]

→ $\text{LF}(6) = 3, L(3) = "n"$ ["nana\$"]

→ $\text{LF}(3) = 7, L(7) = "a"$ ["anana\$"]

Decoding

$\text{rank}_b(L, k) = \# \text{occ of } b \text{ in } L, \text{ up to position } k-1$

$L =$	1	2	3	4	5	6	7	C	1	2	5	6
	a	n	n	b	\$	a	a		\$	a	b	m

first line starting with the letter

Last-to-Front Mapping

$\text{LF}(k) = C[L[k]] + \text{rank}_{L[k]}(L, k)$

\$banana
a\$banan
ana\$ban
anana\$b
banana\$
na\$banana
nana\$ba

→ start with \$ (position 5)

→ $\text{LF}(5) = 1$

→ $L[1] = a$ [current decoding: "a\$"]

→ $\text{LF}(1) = 2, L[2] = n$ ["ba\$"]

→ $\text{LF}(2) = 6, L(6) = "a"$ ["ana\$"]

→ $\text{LF}(6) = 3, L(3) = "n"$ ["nana\$"]

→ $\text{LF}(3) = 7, L(7) = "a"$ ["anana\$"]

→ $\text{LF}(7) = 4, L(4) = "b"$ ["banana\$"]

2. Burrows-Wheeler Transform


\$banana
 a\$banan
 ana\$ban
 anana\$b
 banana\$
 na\$bana
 nana\$ba

What is special about the BTW?

- has many repeating characters (WHY?)
- can be run-length compressed!

Imagine the word “the” appears many times in a text.

he...t
 he...t
 he...t
 he...t (“t”, 18733)
 he...t
 he...t
 he...t



- main motivation
- used in “bzip2” compressor

2. Burrows-Wheeler Transform

\$banana
 a\$banan
 ana\$ban
 anana\$b
 banana\$
 na\$bana
 nana\$ba

What is special about the BTW?

→ efficient backward search!

→ counting #occ's of pattern P in $O(|P| \log |S|)$ time!

Backward Search on BWT

$T = \text{banana}\$$

banana\$	\$banana
\$banana	a\$banan
a\$banan	ana\$ban
na\$bana	anana\$b
ana\$ban	banana\$
nana\$ba	na\$bana
anana\$b	nana\$ba

Burrows-Wheeler Transform L of text T

C	\$	a	b	n
	1	2	5	6

	1	2	3
$P =$	a	n	a
$[sp,ep] =$	[2,4]		

Backward search for Pattern $P[1]..P[m]$

→ Initial range: $[sp,ep]$ with $sp = C[P[m]]$ and $ep = C[P[m]+1]-1$

Then $[s,e]$ with

$$s = C[P[i]] + \text{rank}_{L[i]}(L, sp-1)$$

$$e = C[P[i]] + \text{rank}_{L[i]}(L, ep) - 1$$

Backward Search on BWT

T = banana\$

banana\$
\$banana
a\$banan
na\$bana
ana\$ban
nana\$ba
anana\$b

\$banana
a\$banan
ana\$ban
anana\$b
banana\$
na\$bana
nana\$ba

Burrows-Wheeler Transform L of text T

C \$ a b n
1 2 5 6

P = ana
[sp,ep] = [2,4]

$$s = C["n"] + \text{rank}_n(L, 1) \\ = 6 + 0 = 6$$

$$e = 6 + \text{rank}_n(L, 4) - 1 \\ = 6 + 2 - 1 = 7$$

Backward search for Pattern P[1]..P[m]

$$s = C[P[i]] + \text{rank}_{L[i]}(L, sp-1) \\ e = C[P[i]] + \text{rank}_{L[i]}(L, ep) - 1$$

Backward Search on BWT

T = banana\$

banana\$	\$banana
\$banana	a\$banan
a\$banan	ana\$ban
na\$bana	anana\$b
ana\$ban	banana\$
nana\$ba	na\$bana
anana\$b	nana\$ba

Burrows-Wheeler Transform L of text T

C	\$	a	b	n
	1	2	5	6

P =	1	2	3
	a	n	a

[sp,ep] = [2,4]
sp=6
ep=7

$$s = C["a"] + \text{rank}_a(L, 5) \\ = 2 + 1 = 3$$

$$e = 1 + \text{rank}_a(L, 7) = \\ 2 + 3 - 1 = 4$$

Backward search for Pattern P[1]..P[m]

$$s = C[P[i]] + \text{rank}_{L[i]}(L, sp-1) \\ e = C[P[i]] + \text{rank}_{L[i]}(L, ep) - 1$$

Done!

[3,4]=final range
→ 2 occs of "ana"

BWT Construction

→ use the suffix array $SA(T)$!

→ $BWT[k] = T[SA[k] - 1]$ (assuming $T[0]=\$$)

BWT Construction

→ use the suffix array SA(T)!

→ $BWT[k] = T[SA[k] - 1]$ (assuming $T[0]=\$$)

e.g.

1234567
 SA[banana\$] = [7, 6, 4, 2, 1, 5, 3]

T[6] T[5] T[3] T[1] T[0] T[4] T[2]
 a n n b \$ a a

\$banana
 a\$banan
 ana\$ban
 anana\$b
 banana\$
 na\$bana
 nana\$ba

BWT Construction

- use the suffix array $SA(T)$!
- $BWT[k] = T[SA[k] - 1]$ (assuming $T[0]=\$$)
- explain why this equation is correct!

```
$banana  
a$banan  
ana$ban  
anana$b  
banana$  
na$bana  
nana$ba
```

END

Lecture 16