

Simple Linear-Time Off-Line Text Compression by Longest-First Substitution

Ryosuke Nakamura, Hideo Bannai,
Shunsuke Inenaga, Masayuki Takeda

Kyushu University, Japan

Lossless Compression

Uncompressed



Compressed

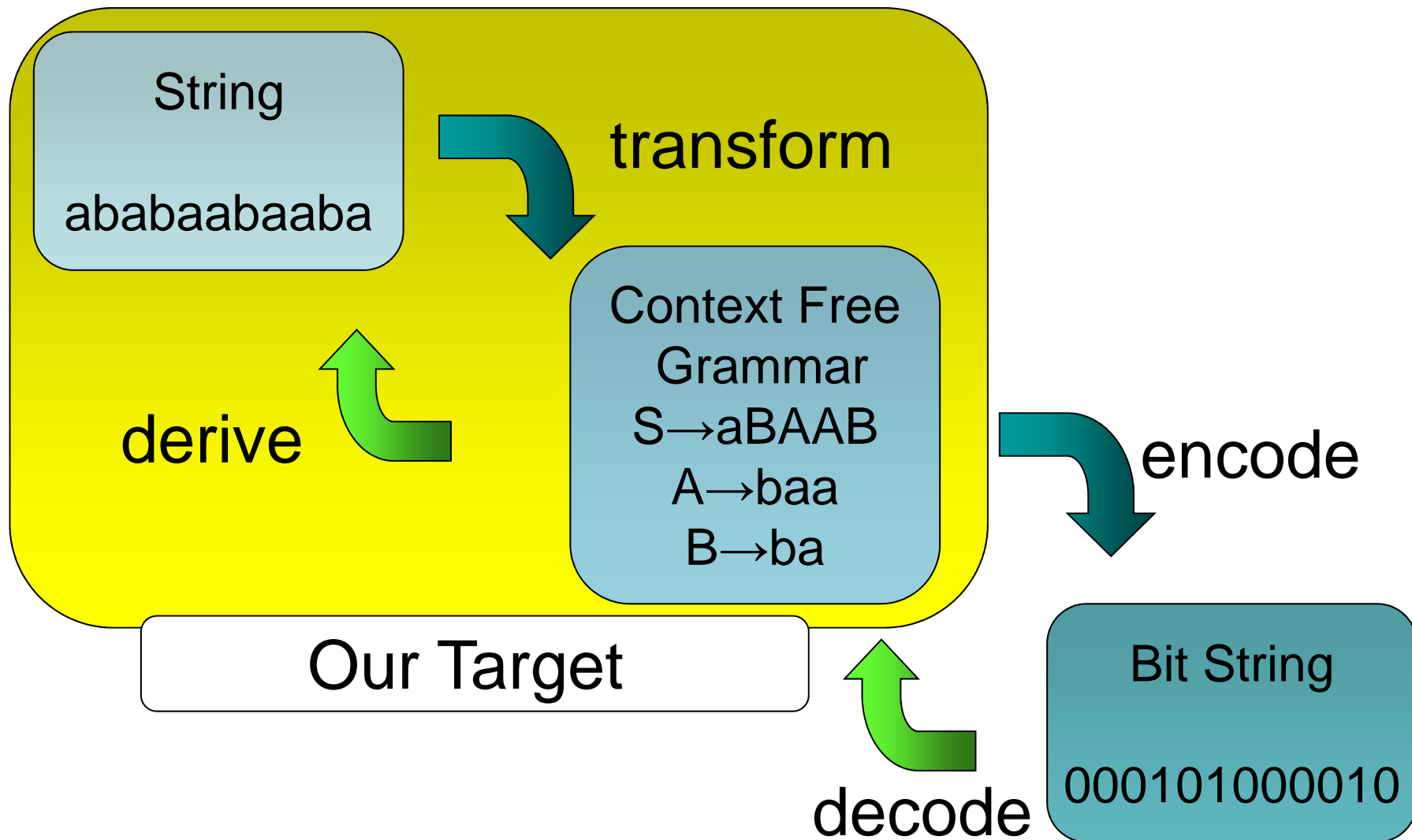


Lossless Data Compression

- Dictionary-based Compression
 - LZ77 [Ziv and Lempel, '77]
- BW Transform-based Compression [Burrows and Wheeler, '94]
 - Bzip
- Grammar-based Compression
 - SEQUITUR [Nevill-Manning and Witten '97]
 - RE-PAIR [Larrson and Moffat '00]

Fast compressed pattern matching [Kida et al. '03]

Grammar-based Compression



Grammar-based Compression

- There exist more than one grammar generating the same string.
- Input string : aaaaaaaaaabaaaaaaaaa

$S \rightarrow \text{aaaaaaaaabaaaaaaaaa}$

$S \rightarrow AbA, A \rightarrow BB, B \rightarrow \text{aaaa}$

$S \rightarrow AAbAA, A \rightarrow \text{aaaa}$

$S \rightarrow AAAAbAAAA, A \rightarrow \text{aa}$

- We want a small grammar.

Grammar-based Compression

- There exist more than one grammar generating the same string
- Input string: $aaabaaa$

Finding the smallest grammar is NP-hard [Storer '77]

$S \rightarrow aaabaaa$

$S \rightarrow AAb$

$S \rightarrow AAAAbAAA, A \rightarrow a$

- We want a small grammar.

Greedy Method

- Recursively replaces frequent and/or long factors by non-terminal symbols.

Input string: **aaaaaaaaabaaaaaaaa**

$\left\{ \begin{array}{l} S \rightarrow XXXXbXXXX \\ X \rightarrow aa \end{array} \right.$

Replaced string : **XXXXXbXXXXX**

$\left\{ \begin{array}{l} S \rightarrow YbY \\ X \rightarrow aa \\ Y \rightarrow XXXX \end{array} \right.$

Greedy Method

- Most Frequent First Substitution (MFFS)
recursively replaces the most frequent factors
by non-terminal symbols



- Longest First Substitution (LFS)
recursively replaces the longest repeating
factors by non-terminal symbols



Greedy Method

- Most Frequent First Substitution (MFFS)

RE-PAIR : linear time algorithm for MFFS
[Larsson and Moffat '00]

- Longest First Substitution (LFS)
recursively replaces the longest repeating factors by non-terminal symbols



Text Compression by Longest First Substitution

[Bentley and McIlroy '99]

- Recursively replaces
Longest Repeating Factors (*LRFs*)
by non-terminal symbols

$S \rightarrow$ acdaabcacabdcaababdcaad

aab ac abdcaa ← longest repeating factor

$S \rightarrow$ acdaabcacXbXd, X \rightarrow abdcaa

$S \rightarrow$ YdaabcYXbXd, X \rightarrow abdcaa, Y \rightarrow ac

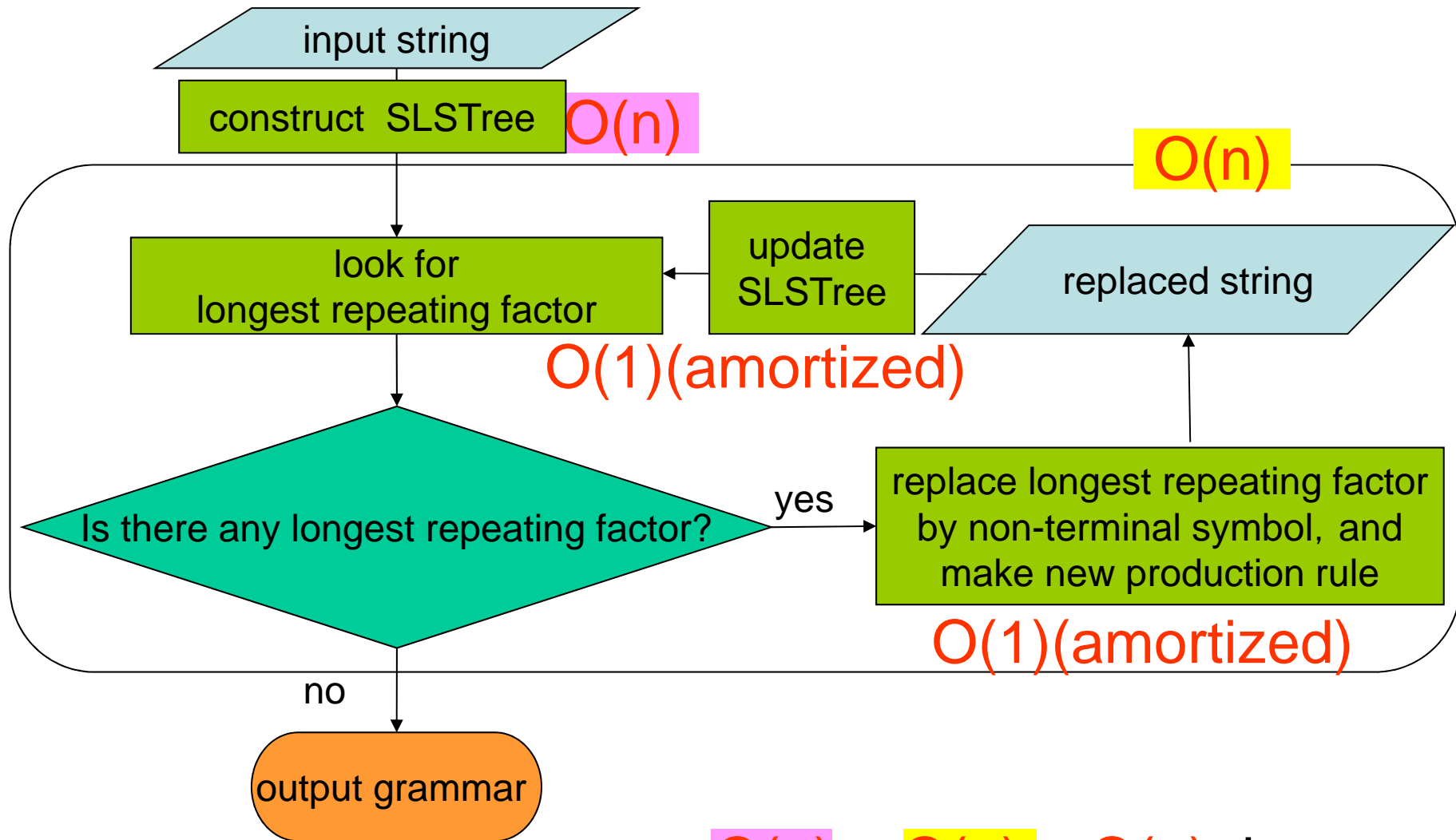
Time Complexities of Longest First Substitution Approach

- n : length of input string
 - naive algorithm : $O(n^4)$ time ☹️
 - algorithm using minimal augmented suffix trees [Brodal et al. '02] : $O(n^2 \log n)$: time 😐
 - our algorithm : $O(n)$ time 😊

Our algorithm

- Uses a new data structure **sparse lazy suffix trees (SLSTrees)** based on suffix trees [Weiner '73];
- Runs in linear time and space.

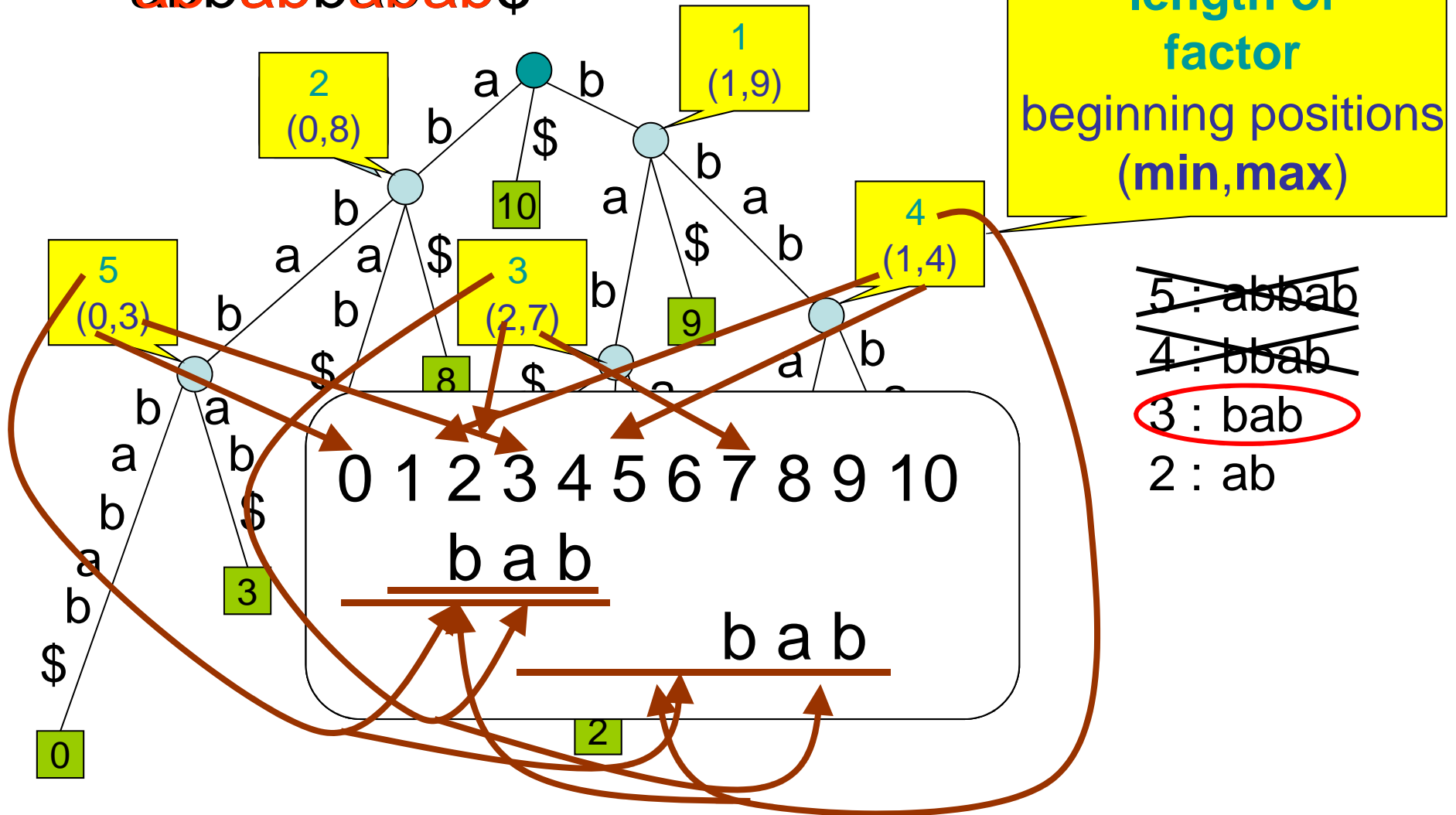
Algorithm Flow Chart



$$O(n) + O(n) = O(n) \text{ time}$$

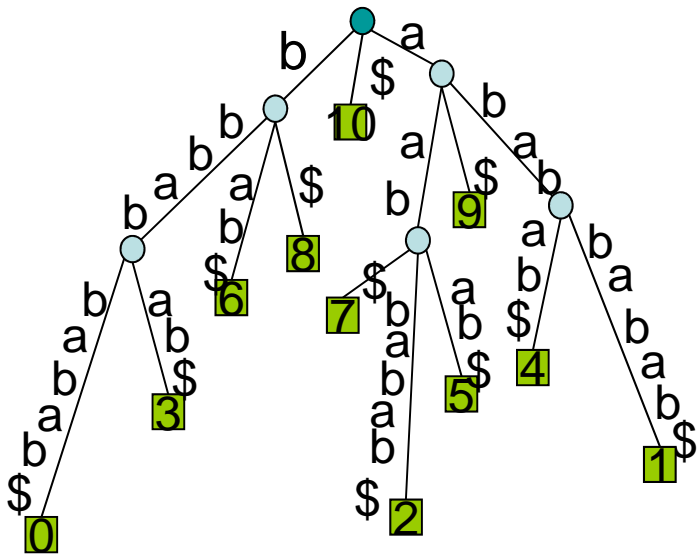
How to find Longest Repeating Factor

012345678910
ababb**ab**ab\$

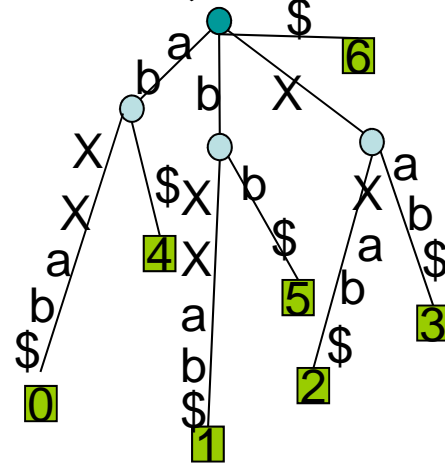


If we use suffix tree...

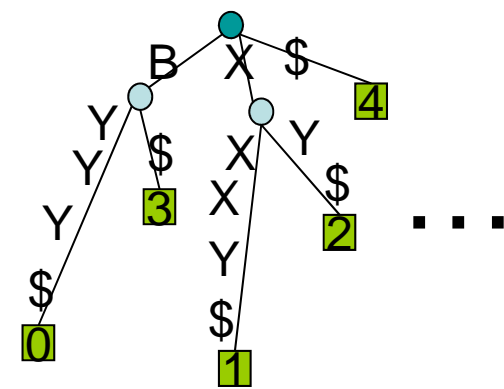
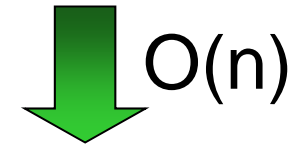
abbabbabab\$



abXXab\$



YXXY\$



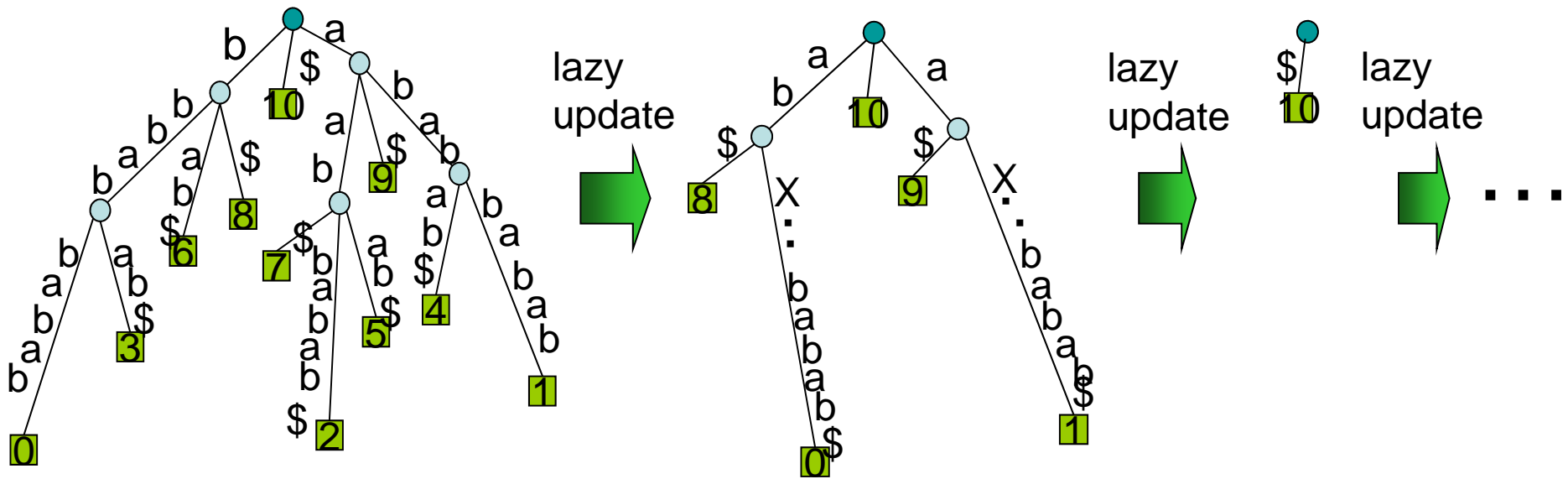
$O(n^2)$

Using sparse lazy suffix tree

abbabbabab\$

abXXab\$

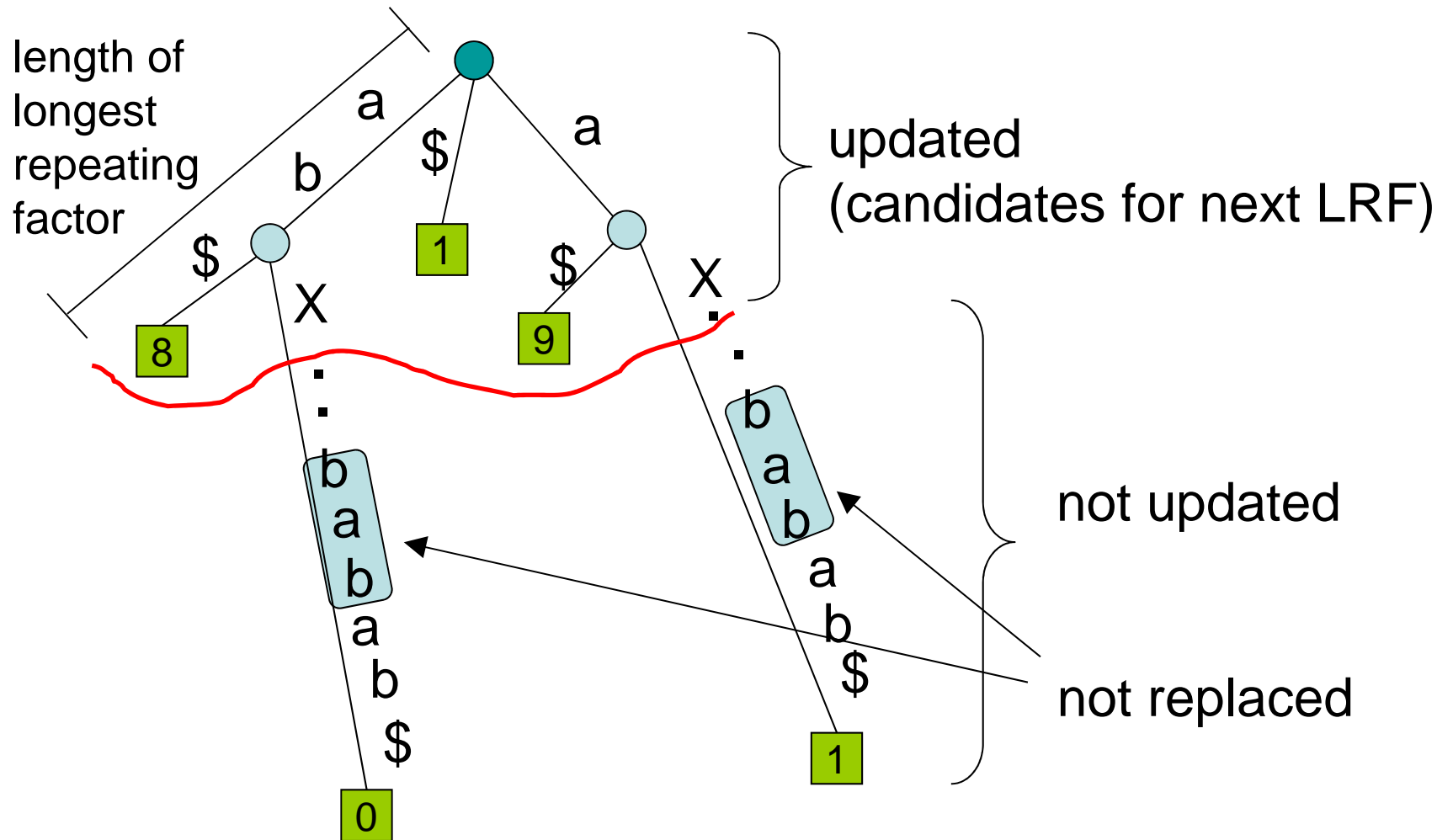
YXXY\$



What is laziness?

012345678910
 abX...X...ab \$

length of next longest repeating factor \leq
 length of current longest repeating factor

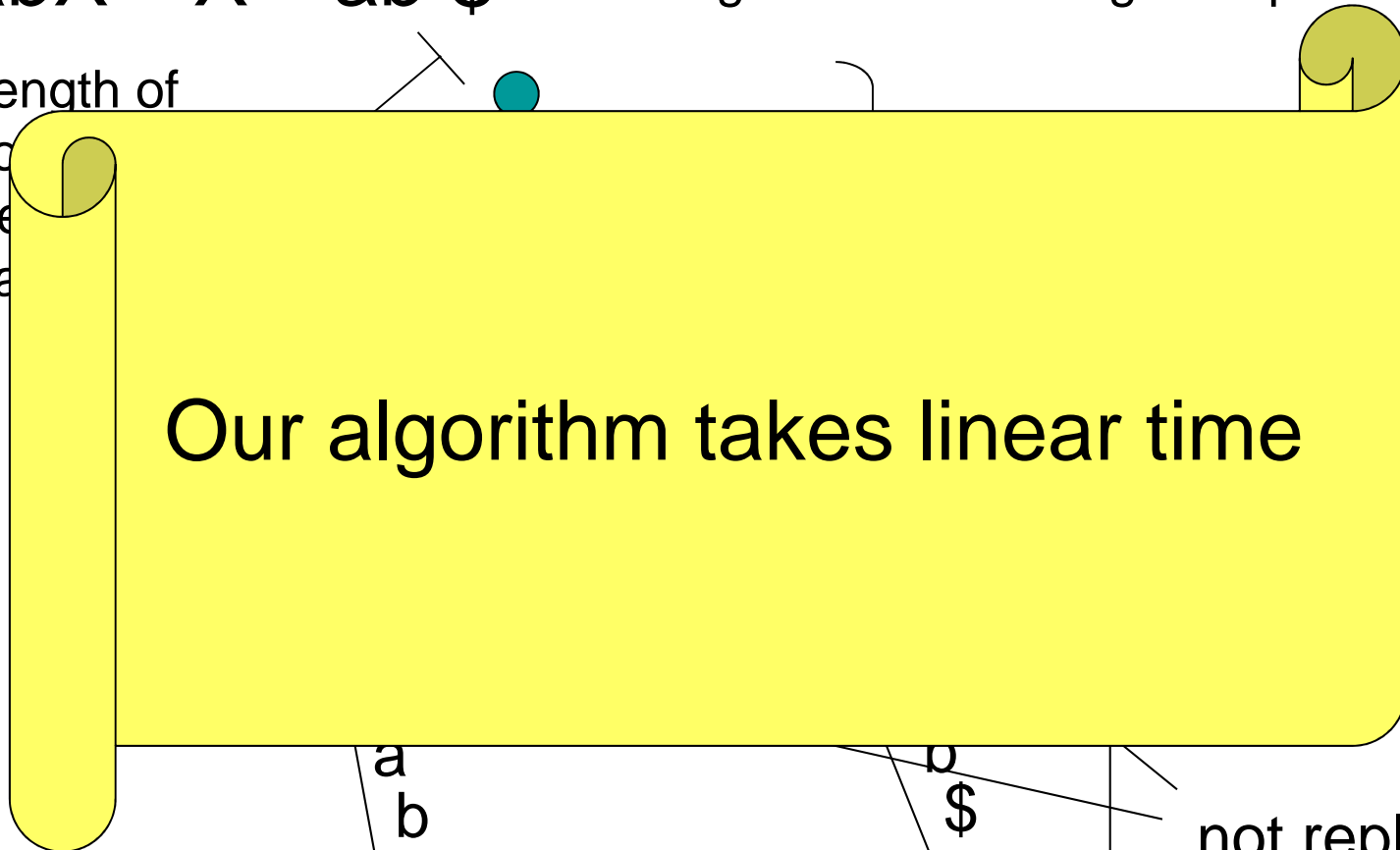


What is laziness?

012345678910
abX...X...ab \$

length of next longest repeating factor \leq
length of current longest repeating factor

length of
lo
re
fa



next LRF)

ated

a
b
\$
0

b
\$
1

not replaced

Smaller grammar with Longest First Substitution (LFS2)

- Consider the right side of production rules constructed by the LFS method.

abbabbabbabbabab\$

$$\left(\begin{array}{l} S \rightarrow XYbYX\$ \\ X \rightarrow a\text{bbab} \\ Y \rightarrow \text{bbab} \end{array} \right) \quad \left(\begin{array}{l} S \rightarrow X\text{bbabbbab}X\$ \\ X \rightarrow \text{abbab} \\ X\text{bbabbbab}X\$_1\text{bbab}\$_2 \end{array} \right)$$

want to replace with Y

$$\left(\begin{array}{l} S \rightarrow XYbYX\$ \\ X \rightarrow aY \\ Y \rightarrow \text{bbab} \end{array} \right)$$

Smaller grammar with Longest First Substitution (LFS2)

- Consider the right side of production rules constructed by the LFS method.

abbabbabbabbabab\$

LFS

$$\begin{cases} S \rightarrow XYbYX\$ \\ X \rightarrow \text{abbab} \\ Y \rightarrow \text{bbab} \end{cases}$$

15

grammar size

LFS2

$$\begin{cases} S \rightarrow XYbYX\$ \\ X \rightarrow aY \\ Y \rightarrow \text{bbab} \end{cases}$$

12

grammar size : total length of the right side of production rules.

Smaller grammar with Longest First Substitution (LFS2)

- Consider the right side of production rules constructed by the LFS method.

abba

We developed a linear
time algorithm for LFS2.

bYX\$

b

grammar size

15

12

grammar size : total length of the right side of production rules.

Comparison of grammar sizes

- Input texts : Canterbury Corpus

File	Size (Bytes)	total grammar size		
		MFFS	LFS	LFS2
alice29.txt	152090	38750	88333	45225
asyoulik.txt	125179	35245	74747	41755
cp.html	24603	8006	14559	7977
fields.c	11150	3535	6525	3307
grammar.lsp	3721	1597	2332	1431
kennedy.xls	1029744	165589	291536	166250
lcet10.txt	426754	84923	235112	103602
plrabn12.txt	481861	116128	276714	144078
ptt5	513216	42813	266040	47885
sum	38240	13023	20846	12103
xargs.1	4227	2096	2772	1906

Conclusions

- We developed
 - a linear time algorithm for longest first substitution (LFS);
 - a linear time algorithm for LFS2.
- LFS2 generates smaller grammars than MFFS for small files.

Future Work

- Efficient encoding of LFS grammars to bit strings.
- Compressed pattern matching specialized for LFS compressed texts.
 - The right side of most production rules consists only of terminal symbols (not LFS2).
e.g., $S \rightarrow YdaabcYXbXd$, $X \rightarrow abdcaa$, $Y \rightarrow ac$
There should be efficient pattern matching algorithms for this.