

This is a slide for introducing my work to my ex-supervisors, which was made on August 15, 2016. This slide implies that core idea in my paper had been developed at this time.

Note that green dialogs are added for notes.

Note that this slide contains English and also Japanese.

In-Place Linear Time Suffix Array Construction

2016/09/15

August 15, 2016

後藤啓介

Keisuke Goto

A topic which is not related to the research

■ 6月 [redacted] とBBQへ

- [redacted] 「いまゴトーさんのLZ77をもっと省領域に計算する方法を研究してみます」
- 後藤「へー(いやー、あれはSAのアルゴリズムを元にしてるし、まずはそっちの問題を解決しなきゃだし、SAのin-place計算なんて無理そうだし、、、あれもしかして解ける?)」
- 後々考えると、[redacted] が言っていたのは線形時間じゃないアルゴリズムだった気がする

This is explaining why I started this research in Japanese

Suffix Array

For a string **T** of length N , the Suffix Array **SA** is an array that stores all suffixes of **T** in lexicographic order.

SA[i] = j indicate the suffix **T** _{j} starting at j is the i -th lexicographically smallest suffix.

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

i	SA [i]	T _{SA[i]}
1	7	\$
2	6	a\$
3	4	ana\$
4	2	anana\$
5	1	banana\$
6	5	na\$
7	3	nana\$

Suffix **T**₂ is 4-th smallest suffix,
and it is stored at **SA**[4]

- Though **SA** only requires $N \log N$ bits, most linear time construction algorithms require the double space of **SA**

Theorem [Nong, 2013]

The suffix array of a string of length N can be computed in linear time and $\sigma \log N + O(\log N)$ bits of working space, where σ is the number of distinct characters in the string

The algorithm requires close to constant space when σ is small, but it still requires close to the double space of **SA** when σ is large

■ 昨日軽く調べたら色々見つかりました(汗

■ [Rahmann and Smula, ???] バイオ系の人？出典分からず

- Simple In-Place Suffix Array Construction by Walking along Burrows-Wheeler Transforms
- BWTがあればSAとLCPを線形時間in-placeに計算できるよ(たぶんgeneral alphabet)

■ [Franceschini1 and Muthukrishnan, 2007]

- General alphabetについて $O(N \log N)$ 時間でin-placeにSAを計算できるよ

■ [Nong, 2013]

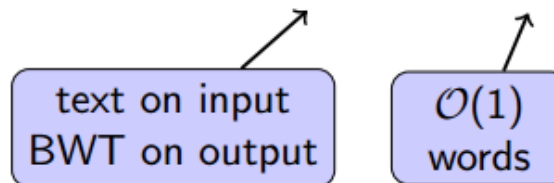
- Constant alphabetについて線形時間でin-placeにSAを計算できるよ

■ [Chrochemore+, 2013]

- $O(n^2)$ 時間でin-placeにBWTが計算できるよ

Computing BWT

	time	space (bits)
Suffix array construction		
many algorithms	$\mathcal{O}(n)$	$> n \log n + n \log \sigma$
Franceschini & Muthu, 2007	$\mathcal{O}(n \log n)$	$n \log n + n \log \sigma + \mathcal{O}(\log n)$
Direct BWT construction		
K, 2007	$\mathcal{O}(n/\epsilon^2)$	$2n \log \sigma + \mathcal{O}(\epsilon n \log n)$
Okanohara & Sadakane, 2009	$\mathcal{O}(n)$	$\mathcal{O}(n \log \sigma \log \log_{\sigma} n)$
Succinct index construction		
Hon & al., 2007	$\mathcal{O}(n \log n)$	$n \log \sigma + \mathcal{O}(n H_0)$
Hon & al., 2009	$\mathcal{O}(n \log \log \sigma)$	$\mathcal{O}(n \log \sigma)$
In-place BWT construction		
This talk	$\mathcal{O}(n^2)$	$n \log \sigma + \mathcal{O}(\log n)$



- We propose the first algorithm to construct the suffix array in linear time and in-place
 - Our algorithm is based on Nong's algorithm
 - Our strategy is first construct some suffix arrays for subset of suffixes, and then merge them by any stable in-place merge algorithm for two sorted arrays

Preliminaries of Nong's Algorithm

- We call a suffix T_i is **small type suffix (S-suffix)** if $i = N$ or $T_i < T_{i+1}$, and **large type suffix (L-suffix)** otherwise
- We call T_i also **left most S-suffix (LMS-suffix)** if T_i is S-suffix and $i=1$ or T_{i-1} is L-suffix, and similarly call T_i **left most L-suffix (LML-suffix)** if T_i is L-suffix

1	2	3	4	5	6	7
<u>b</u>	<u>a</u>	<u>n</u>	<u>a</u>	<u>n</u>	a	\$

$T[N]=\$$ is the most smallest character and does not appear in $T[1..N-1]$

Property

- If $T[i] \neq T[i+1]$, the type of T_i equals to T_{i+1}
- By a right-to-left scan on T , each type of $T[i]$ can be obtained

Preliminaries of Nong's Algorithm

- We call a suffix T_i is **small type suffix (S-suffix)** if $i = N$ or $T_i < T_{i+1}$, and **large type suffix (L-suffix)** otherwise
- We call T_i also **left most S-suffix (LMS-suffix)** if T_i is S-suffix and $i=1$ or T_{i-1} is L-suffix, and similarly call T_i **left most L-suffix (LML-suffix)** if T_i is L-suffix

For a character c , we call the interval of suffixes starting with c as c -interval

$T =$

1	2	3	4	5	6	7
<u>b</u>	<u>a</u>	<u>n</u>	<u>a</u>	<u>n</u>	a	\$

n-interval

a-interval

b-interval

n-interval

i	$SA[i]$	$T_{SA[i]}$
1	7	\$
2	6	a\$
3	4	<u>a</u> na\$
4	2	<u>a</u> n <u>a</u> na\$
5	1	<u>b</u> <u>a</u> n <u>a</u> na\$
6	5	<u>n</u> a\$
7	3	<u>n</u> <u>a</u> na\$

L-suffixes must be smaller than S-suffixes in each interval

Overview of Nong's Algorithm

■ Sort LMS-suffixes

Combination of sorting L, S-suffixes

■ Sort LMS-substring

■ If they are not unique

Run in-place

- Make \mathbf{T}' from \mathbf{T} by replacing all LMS-substring with their ranks
- Constructs \mathbf{SA}' of \mathbf{T}' recursively

■ Sort L-suffixes

Use $\sigma \log N + O(\log N)$ bits of space

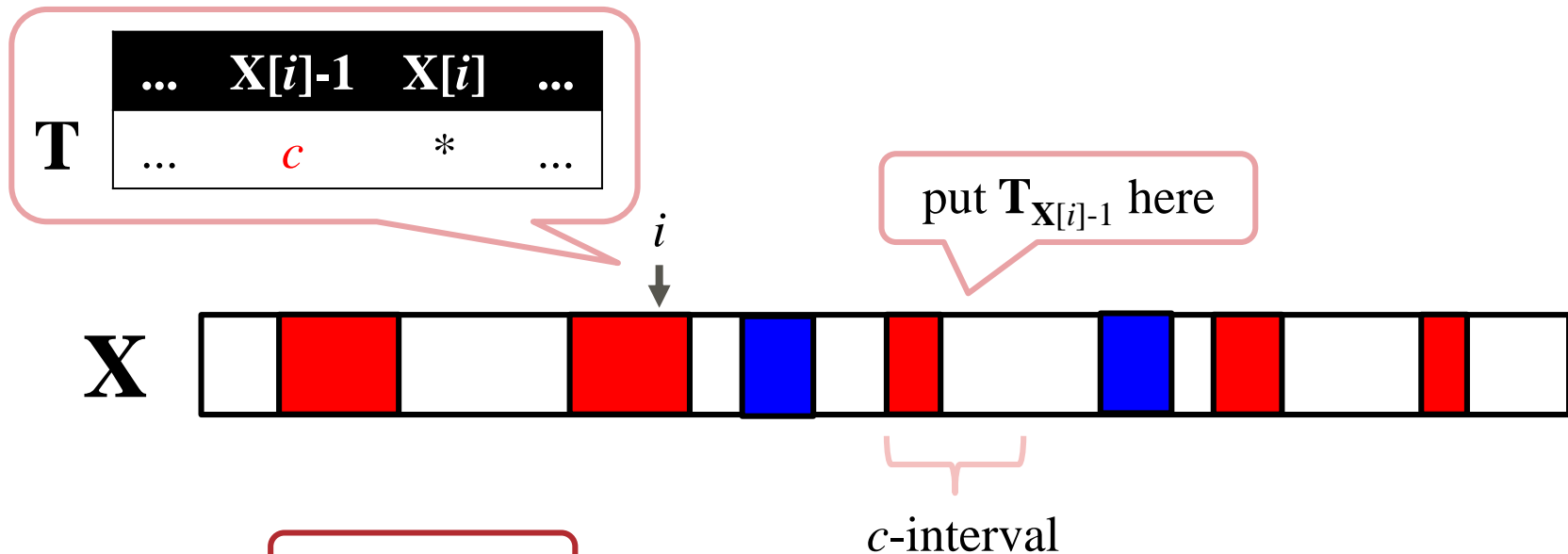
■ Sort S-suffixes

Almost same of sorting L-suffixes

Sorting L-suffixes is the core of algorithm, and also the most space bottleneck part

Sketch of Sorting L-suffixes

- Initially sorted LMS-suffixes are stored in the tail of each c -interval
- By a left-to-right scan on \mathbf{X} ,
 - If $\mathbf{X}[i]$ is not EMPTY and $\mathbf{T}_{\mathbf{X}[i]-1}$ starting with c is L-suffix, put $\mathbf{T}_{\mathbf{X}[i]-1}$ to the left-most-empty position of c -interval on \mathbf{SA}



Correctness

For a L-suffix \mathbf{T}_j stored before $\mathbf{T}_{\mathbf{X}[i]-1}$, $\mathbf{T}_j < \mathbf{T}_{\mathbf{X}[i]-1}$ since \mathbf{T}_j must be put before accessing i , and so $\mathbf{T}_{j+1} < \mathbf{T}_{\mathbf{X}[i]}$

- Initially sorted LMS-suffixes are stored in the tail of each c -interval
 - By a left-to-right scan on \mathbf{X} ,
 - If $\mathbf{X}[i]$ is not EMPTY and $\mathbf{T}_{\mathbf{X}[i]-1}$ starting with c is L-suffix,
put $\mathbf{T}_{\mathbf{X}[i]-1}$ to the left-most-empty position of c -interval on \mathbf{SA}
- Judge by $\mathbf{bkt}[c]$
- Mange by $\mathbf{bkt}[c]$ in $\sigma \log N$ bits

Our algorithm conceptually run in the same way,
but we store $\mathbf{SA}_{s(L)}$ in continuous space of \mathbf{X} ,
and use the remained continuous empty space in \mathbf{X} as \mathbf{bkt}

Notation

- $s(L)$: set of all L-suffixes
- $s(S)$: set of all S-suffixes
- SA_M : suffix array of any set of suffixes M

We firstly explain in the case $\sigma \leq N/2$ and $|s(L)| \leq |s(S)|$, and other case later

mean also $|s(L)| \leq N/2$
since $|s(L)| + |s(S)| = N$

Sketch of Our Algorithm

- $s(Lx')$: the set of largest suffixes starting with c for all characters c
 $s(Lx) = s(L) - s(Lx')$
- $s(LMSx')$: the set of smallest suffixes starting with c such that there are **no** L-suffixes starting with c
 $s(LMSx) = s(LMS) - s(LMSx')$

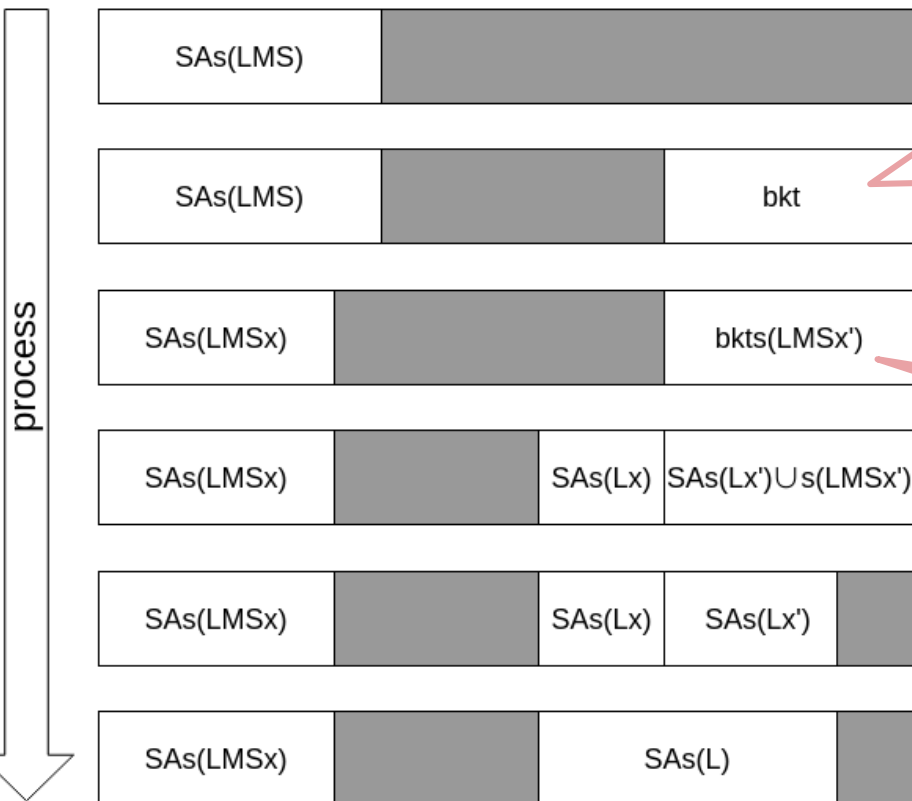
Note that $|s(Lx')| + |s(LMSx')| = \sigma$

サイズと状態だけに注目

bkt[c]: left most empty position of c -interval in $\mathbf{SA}_{s(Lx)}$ if the corresponding L-suffix in $s(LMSx')$ exists, EMPTY

~~Otherwise~~ $|\mathbf{SA}_{s(LMS)}| + |\mathbf{bkt}| \leq N$ since $|s(LMS)| \leq N/2$ and $\sigma \leq N/2$

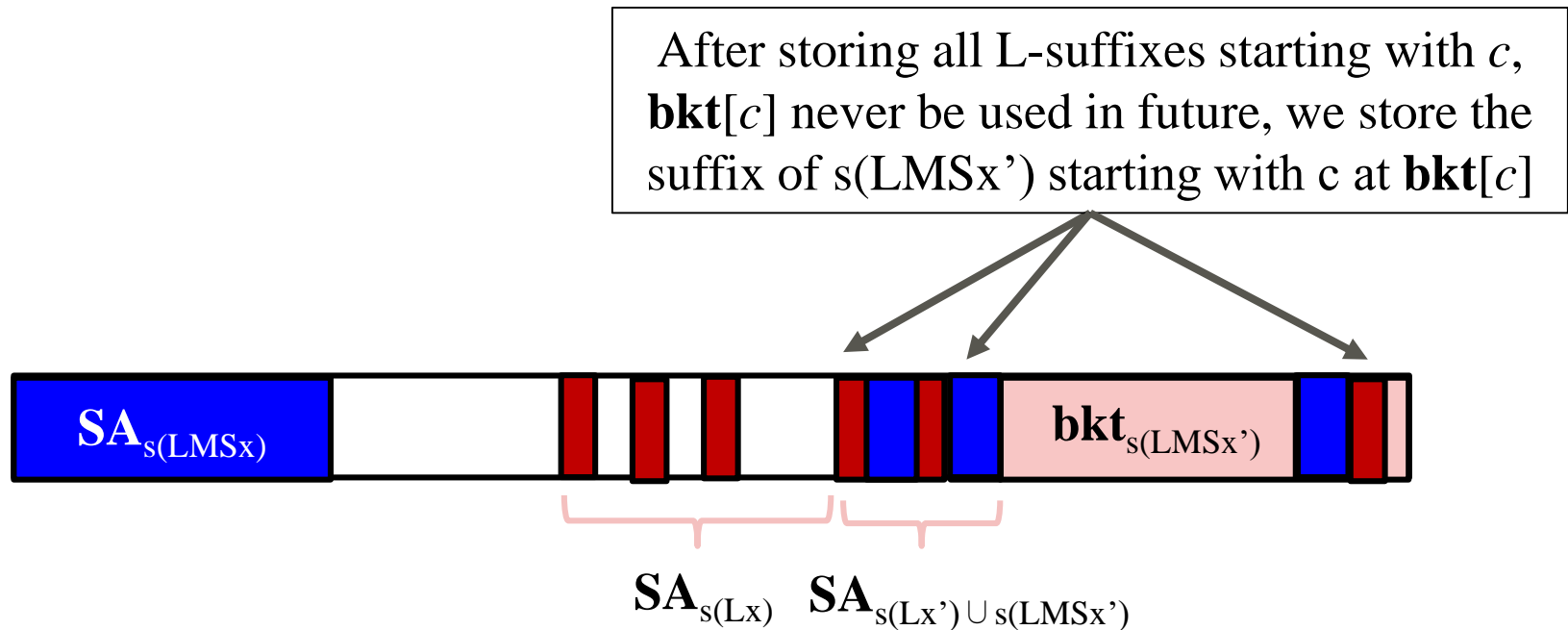
bkt _{$s(LMSx')$} [c]: **bkt**[c] if **bkt**[c] \neq EMPTY, a suffix of $s(LMSx')$ starting with c otherwise



-
- The diagram illustrates the LMS algorithm for a 10-tap FIR filter. The top part shows the initial state with a 10-tap filter. The bottom part shows the state after one iteration, where the filter coefficients are updated based on the error signal.
- Top Diagram (Initial State):**
- The filter has 10 taps, represented by a horizontal bar divided into 10 segments.
 - The first 5 taps are labeled $SA_{s(LMS)}$ (blue).
 - The last 5 taps are labeled bkt (red).
 - The output of the filter is labeled $EMPTY$.
 - Gray arrows indicate the input signal and the output signal.
 - Blue triangles mark the boundaries between the two filter sections.
- Bottom Diagram (State after one iteration):**
- The filter has 10 taps, represented by a horizontal bar divided into 10 segments.
 - The first 5 taps are labeled $s(LMSx)$ (blue).
 - The last 5 taps are labeled $bkt_{s(LMSx')}$ (red).
 - The output of the filter is labeled $s(LMSx')$.
 - Gray arrows indicate the input signal and the output signal.
 - Blue triangles mark the boundaries between the two filter sections.

How to Transit In-Place?

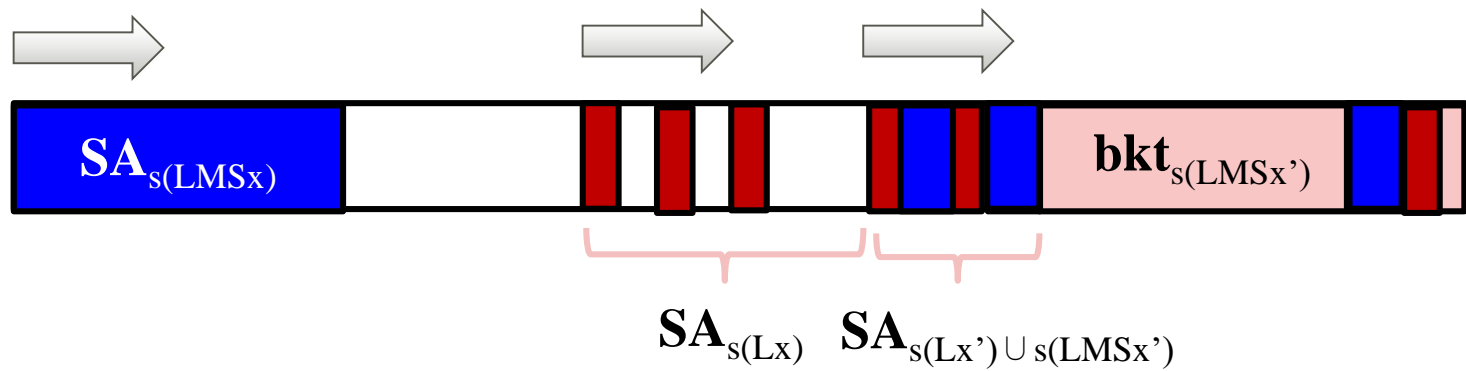
■ $SA_{s(Lx)}$ and $SA_{s(Lx) \cup s(LMSx')}$:



How to Transit In-Place?

$SA_{s(L_X)}$ and $SA_{s(L_X) \cup s(LMS_X')}$:

- By a left-to-right scan on $SA_{s(L_X)}$, $SA_{s(L_X') \cup s(LMS_X')}$, $SA_{s(LMS_X)}$, access all LMS, L-suffixes T_i starting with c lexicographically
 - Judge whether T_{i-1} is L-suffix, if not do nothing
 - We try to put L-suffix T_{i-1} starting with c to $SA_{s(L_X)}[bkt[c]]$

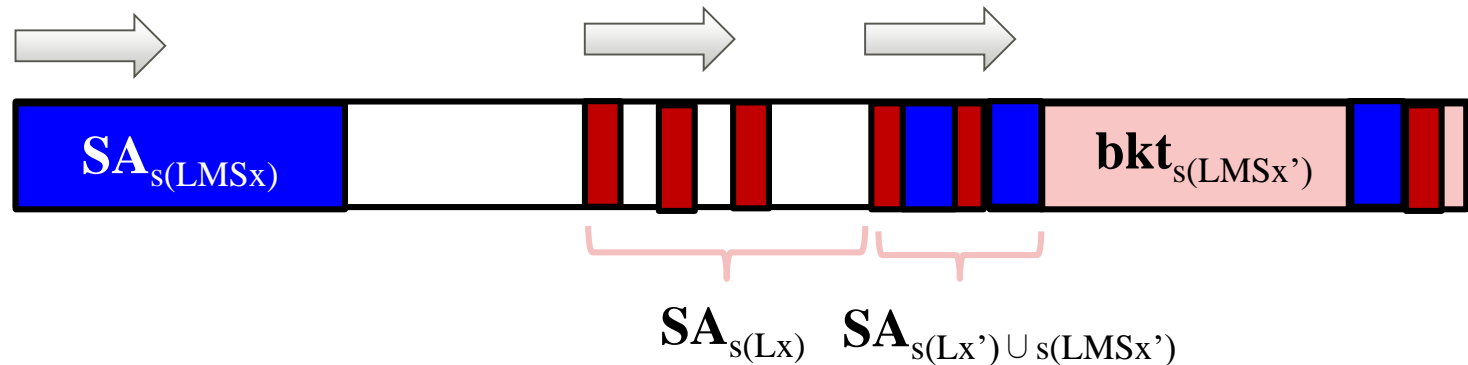


How to Transit In-Place?

$\mathbf{SA}_{s(Lx)}$ and $\mathbf{SA}_{s(Lx) \cup s(LMSx')}$:

Judge whether \mathbf{T}_{i-1} is L-suffix, if not do nothing

- If \mathbf{T}_i is of $\mathbf{SA}_{s(Lx)}$ or $\mathbf{SA}_{s(LMSx)}$, we can determine the type of \mathbf{T}_{i-1} by comparing heading characters of \mathbf{c} of \mathbf{T}_i
- If \mathbf{T}_i is of $\mathbf{SA}_{s(Lx') \cup s(LMSx')}$, the heading characters of \mathbf{T}_{i-1} of \mathbf{T}_i must be different, so we can determine the type of \mathbf{T}_{i-1} by comparing heading characters



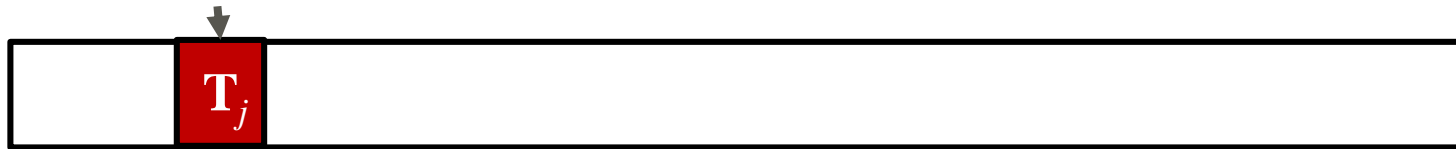
How to Transit In-Place?

$\mathbf{SA}_{s(Lx)}$ and $\mathbf{SA}_{s(Lx) \cup s(LMSx')}$:

We try to put L-suffix \mathbf{T}_{i-1} starting with c to $\mathbf{SA}_{s(Lx)}[\mathbf{bkt}[c]]$

- If $\mathbf{SA}_{s(Lx)}[\mathbf{bkt}[c]]$ is EMPTY, we put \mathbf{T}_{i-1} there
- Otherwise

$\mathbf{SA}[\mathbf{bkt}[c]]$
cannot be stored \mathbf{T}_{i-1}



At least \mathbf{T}_{i-1} or \mathbf{T}_j ($j = \mathbf{SA}_{s(Lx)}[\mathbf{bkt}[c]]$) is of $s(Lx')$, and heading characters of them must be different since $\mathbf{bkt}[c]$ must indicate different positions for each L-suffix starting with c

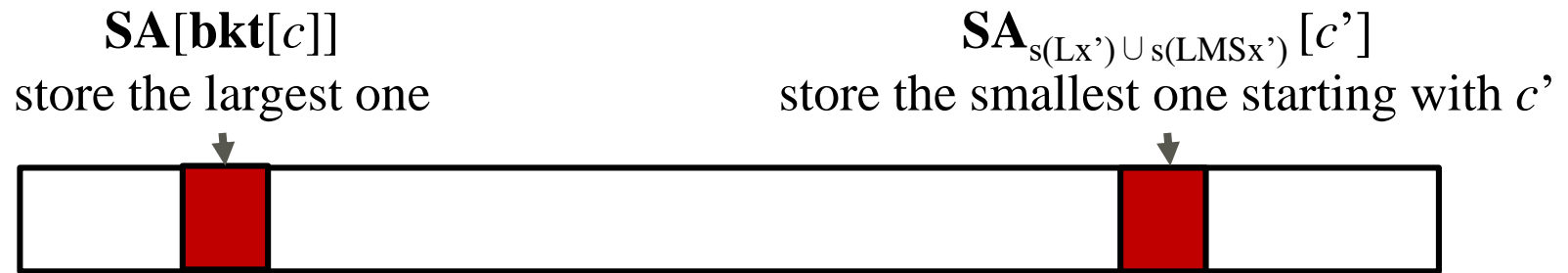
Conflict occurs between a largest suffix starting with $c1$ and a smallest suffix starting with $c2 > c1$

How to Transit In-Place?

$\mathbf{SA}_{s(L_X)}$ and $\mathbf{SA}_{s(L_X) \cup s(LMS_X')}$:

We try to put L-suffix \mathbf{T}_{i-1} starting with c to $\mathbf{SA}_{s(L_X)}[\mathbf{bkt}[c]]$

- If $\mathbf{SA}_{s(L_X)}[\mathbf{bkt}[c]]$ is EMPTY, we put \mathbf{T}_{i-1} there
- Otherwise, we put the smallest of $\mathbf{T}(\mathbf{SA}_{s(L_X)}[\mathbf{bkt}[c]])$ and \mathbf{T}_i to $\mathbf{SA}_{s(L_X') \cup s(LMS_X')}[c']$ and the other into $\mathbf{SA}_{s(L_X)}[\mathbf{bkt}[c]]$, where c' is the starting character of smallest one

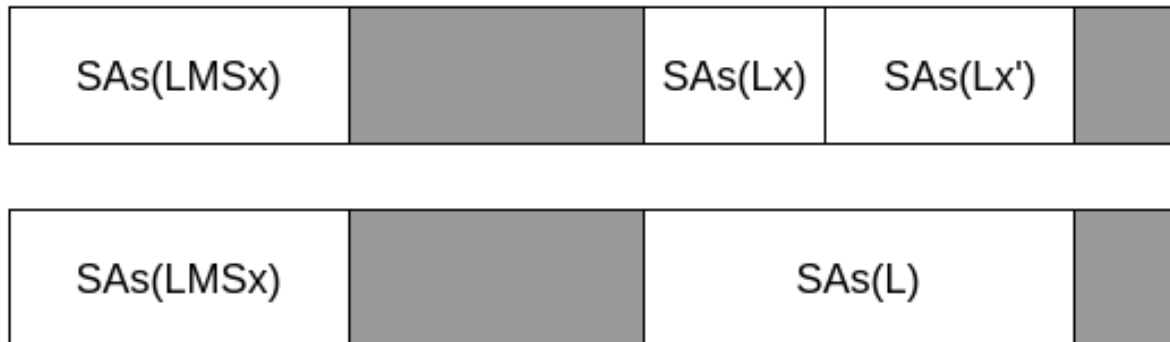


Merge Two Suffix Arrays

$\mathbf{SA}_{s(L)}$:

Observation: In each c -interval in $\mathbf{SA}_{s(L)}$, Each L-suffix of $\mathbf{SA}_{s(Lx')}$ must be stored at the right side of all L-suffixes of $\mathbf{SA}_{s(Lx)}$

- We can make $\mathbf{SA}_{s(L)}$ by a stable merge of $\mathbf{SA}_{s(Lx)}$ and $\mathbf{SA}_{s(Lx')}$ with $\mathbf{T}[\mathbf{SA}_{s(Lx)}[i]]$ and $\mathbf{T}[\mathbf{SA}_{s(Lx')}[i]]$ as keys



Merge Two Suffix Arrays

$\mathbf{SA}_{s(L)} :$

Observation: In each c -interval in $\mathbf{SA}_{s(L)}$, Each L -suffix of $\mathbf{SA}_{s(Lx')}$ must be stored at the right side of all L -suffixes of $\mathbf{SA}_{s(Lx)}$

- We can make $\mathbf{SA}_{s(L)}$ by a stable merge of $\mathbf{SA}_{s(Lx)}$ and $\mathbf{SA}_{s(Lx')}$ with $\mathbf{T}[\mathbf{SA}_{s(Lx)}[i]]$ and $\mathbf{T}[\mathbf{SA}_{s(Lx')}[i]]$ as keys

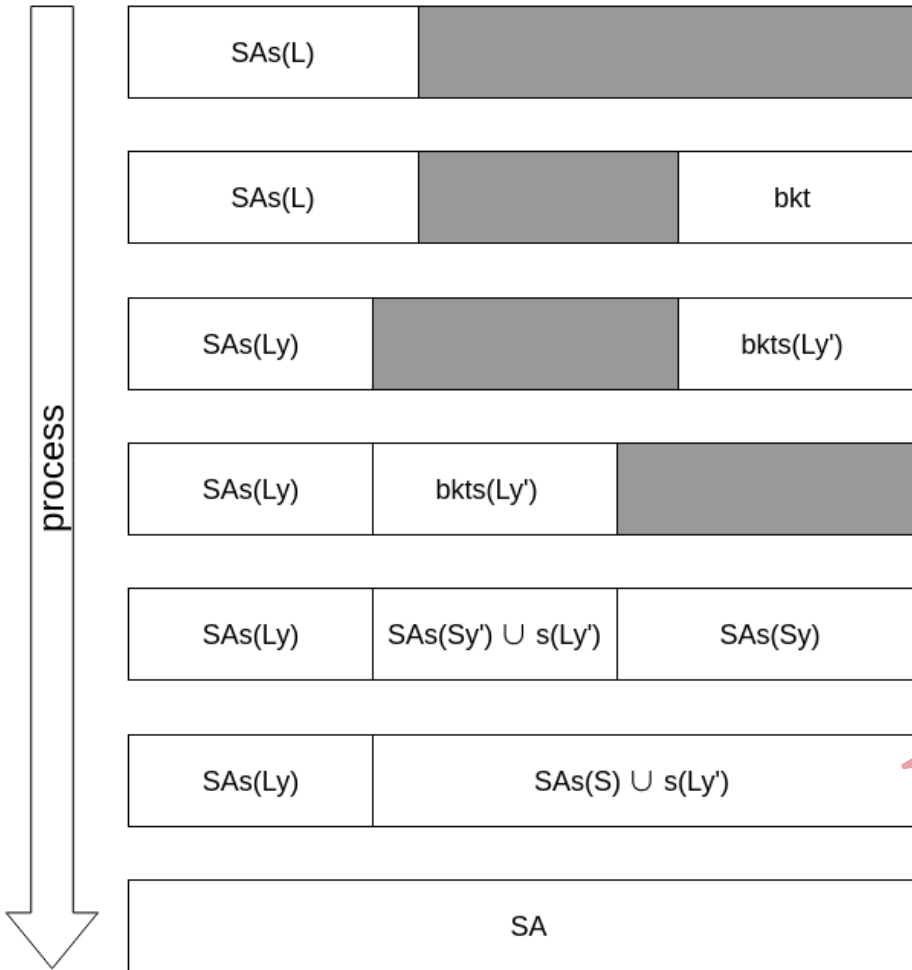
Theorem [Huang and Langston, 1998]

For two integer arrays \mathbf{Y}_1 of length N_1 and \mathbf{Y}_2 of length N_2 , which are stored in \mathbf{Y} of length $N_1 + N_2$ in this order, there is an algorithm to sort the values in \mathbf{Y} stably in $O(N_1 + N_2)$ and in-place

Lemma

There is an algorithm to compute $\mathbf{SA}_{s(L)}$ linear time and in-place when $\sigma \leq N/2$

■ Almost same of sorting L-suffixes

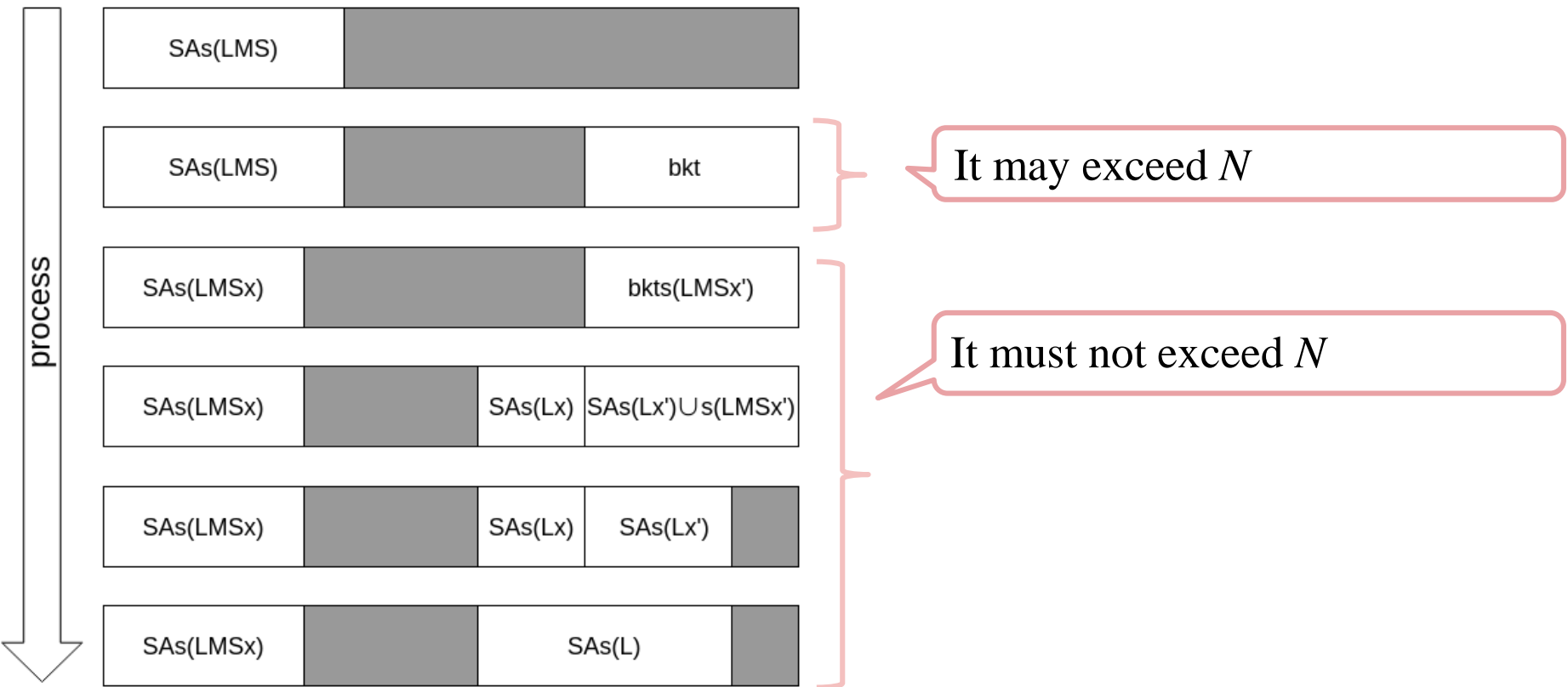


Note that $|\mathbf{SAs}(L)| + |\mathbf{bkt}| \leq N$
 since we assumed $|s(L)| \leq N/2$ and
 $\sigma \leq N/2$

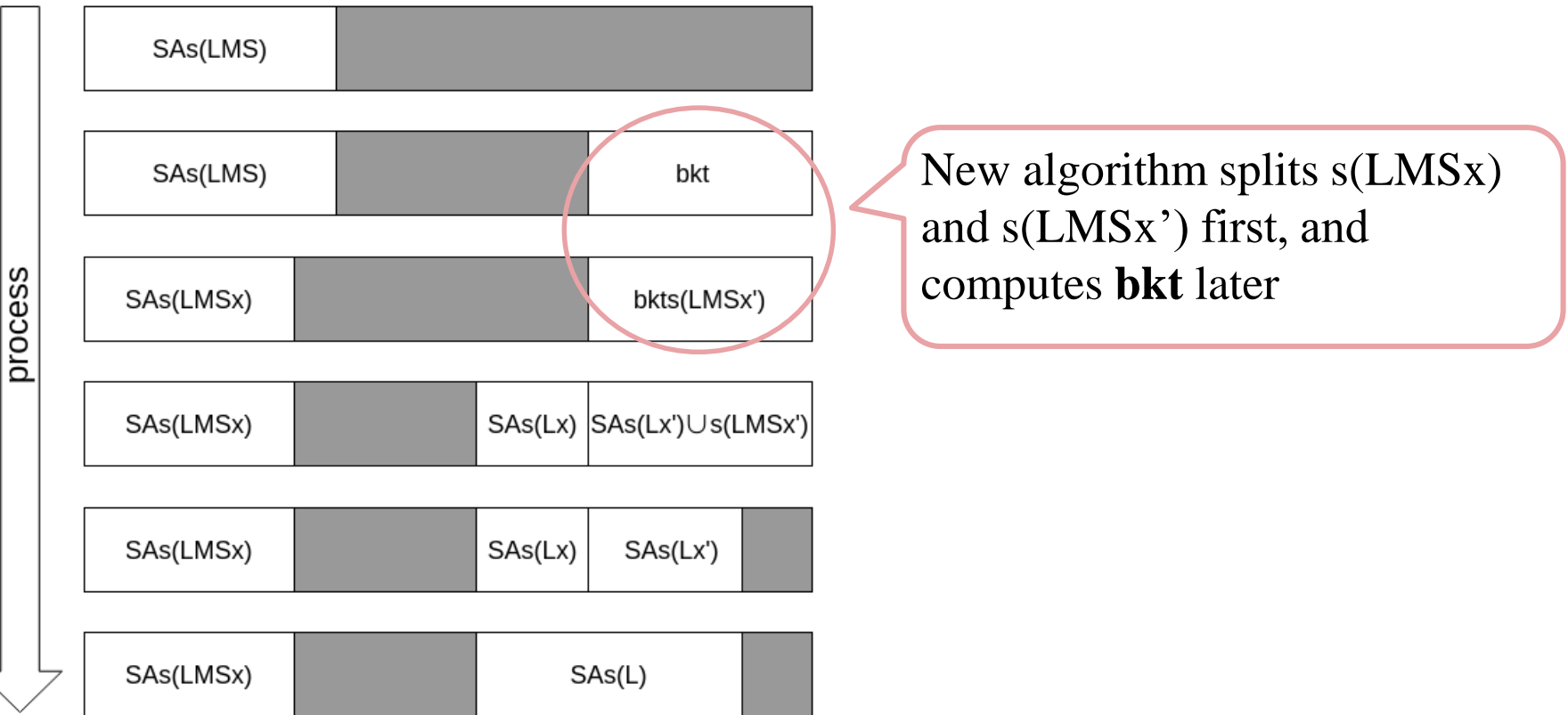
Make **SA** rather than $\mathbf{SA}_{s(S)}$

- If $\sigma > N/2$ or $|s(L)| > |s(S)|$, there may not be sufficient space to store $\mathbf{SA}_{s(LMS)}/\mathbf{SA}_{s(L)}$ and \mathbf{bkt}
- Previous assumption:
 - $\sigma > N/2$ and $|s(L)| \leq |s(S)|$
- We next consider the case $\sigma > N/2$ and $|s(L)| \leq |s(S)|$

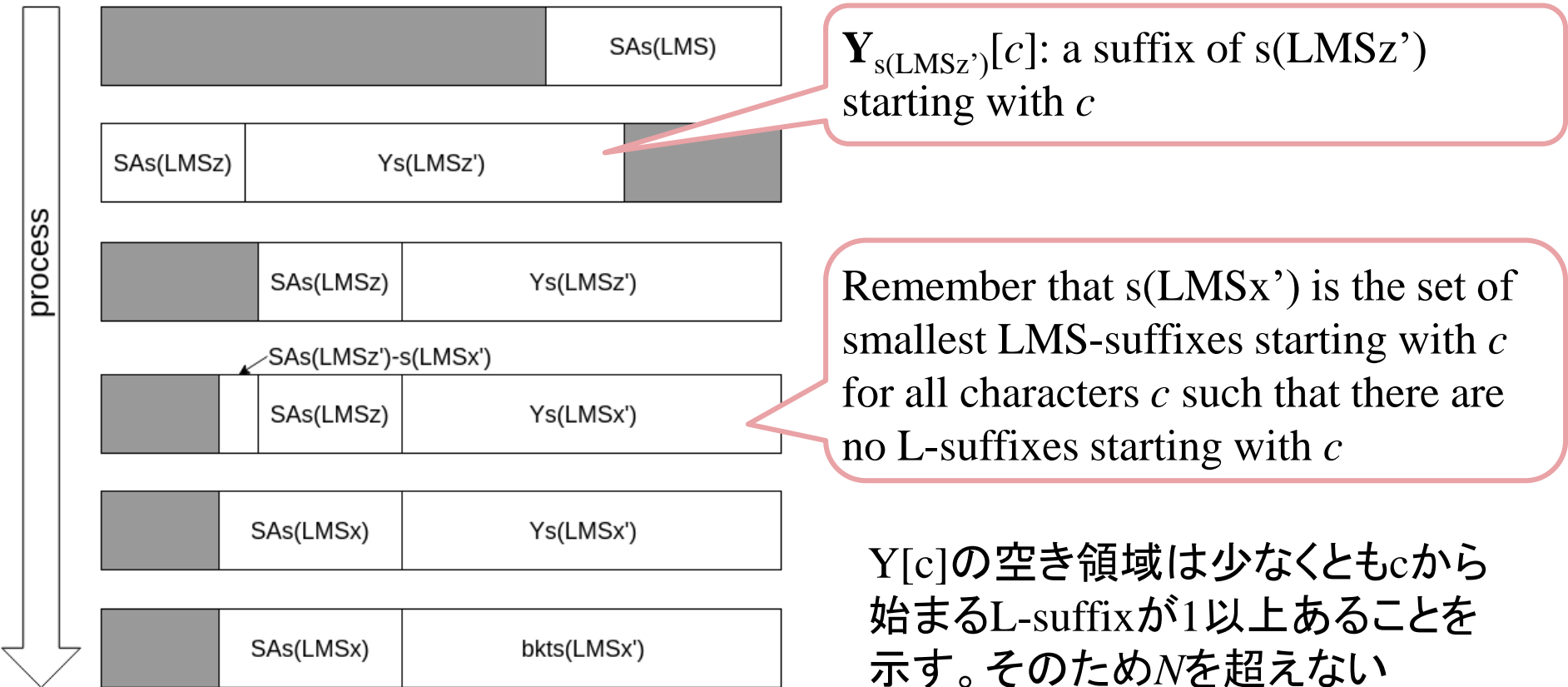
- We can compute $\mathbf{SA}_{s(L)}$ in-place if we can compute $\mathbf{SA}_{s(LMSx)}$ and $\mathbf{bkt}_{s(LMSx')}$ in-place



- We can compute $\mathbf{SA}_{s(L)}$ in-place if we can compute $\mathbf{SA}_{s(LMS_x)}$ and $\mathbf{bkt}_{s(LMS_x')}$ in-place



- $s(\text{LMSz}')$: the set of smallest suffixes starting with c for all characters c
 $s(\text{LMSz}) = s(\text{LMS}) - s(\text{LMSz}')$



I used this notation as the ceil function

Lemma

If $\mathbf{T}[i] \neq \mathbf{T}[i-1]$, i can be represented by $\lceil \log N \rceil - 1$ bits, and back again if we know $\mathbf{T}[i]$ and \mathbf{T}

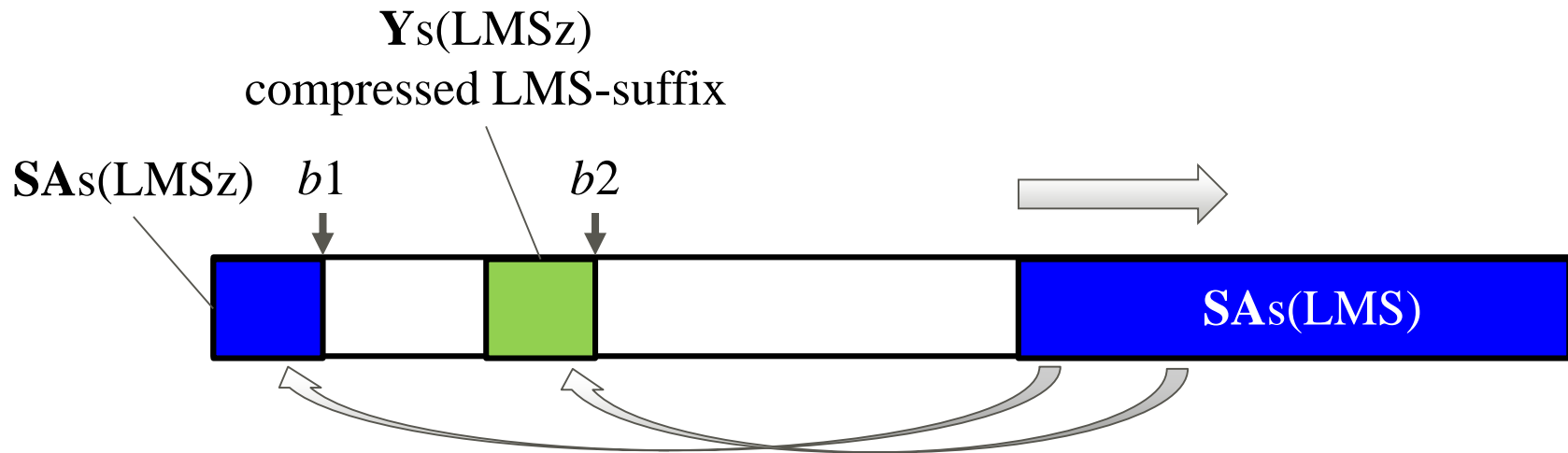
■ Proof

- Let $j = \lceil i/2 \rceil$ and $c = \mathbf{T}[i]$
Note that i is $2j$ or $2j+1$
- If $\mathbf{T}[2j] \neq \mathbf{T}[2j+1]$, the position of either which equals to $\mathbf{T}[i]$ is i
- Otherwise $i = 2j$ since if not, it contradicts the definition

How to Transit In-Place?

■ $\mathbf{SA}_{s(\text{LMSz})}$ and $\mathbf{Y}_{s(\text{LMSz})}$:

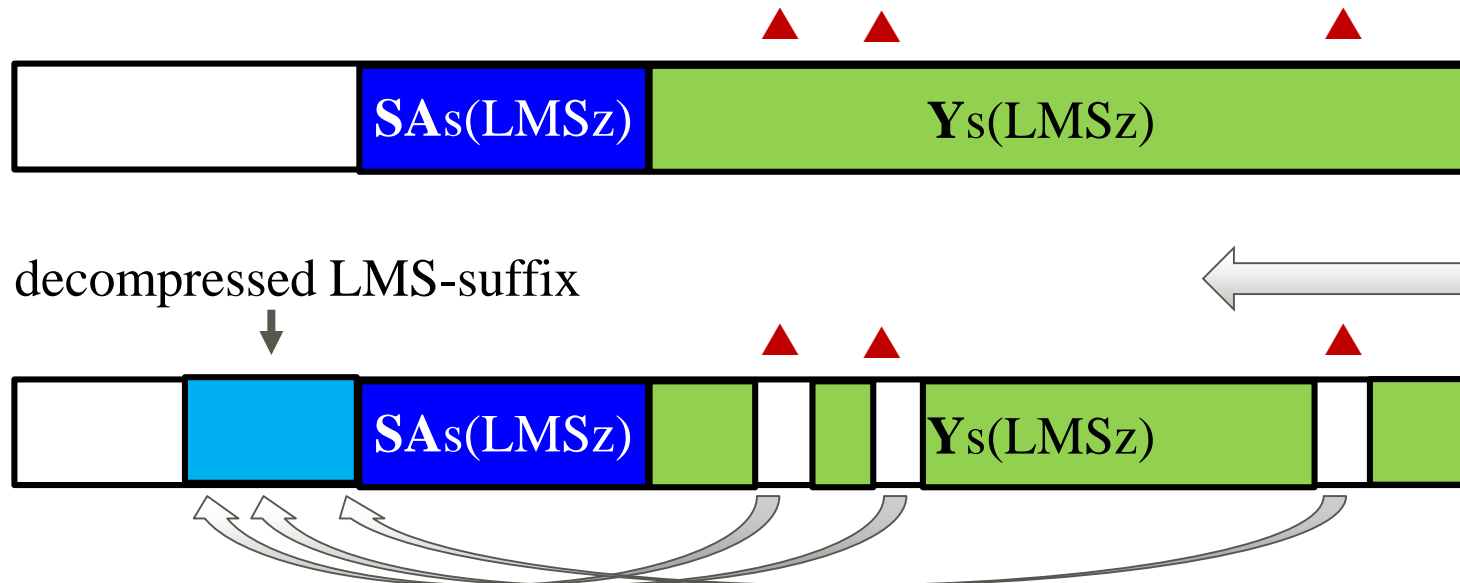
- Let $b1=0$, $b2=|s(\text{LMSz})|$
- By a left-to-right scan on $\mathbf{SA}_{s(\text{LMS})}$, when accessing $\mathbf{SA}_{s(\text{LMS})}[i]=j$
 - If $\mathbf{T}_j \in s(\text{LMSz})$, move j to $\mathbf{X}[b1]$, and update $b1+=1$
 - otherwise, move $\lceil j/2 \rceil$ to $\mathbf{X}[b2]$ and update $b2+=1$



How to Transit In-Place?

■ $SA_{s(LMSz')-s(LMSx)}$ and $Y_{s(LMSx)}$:

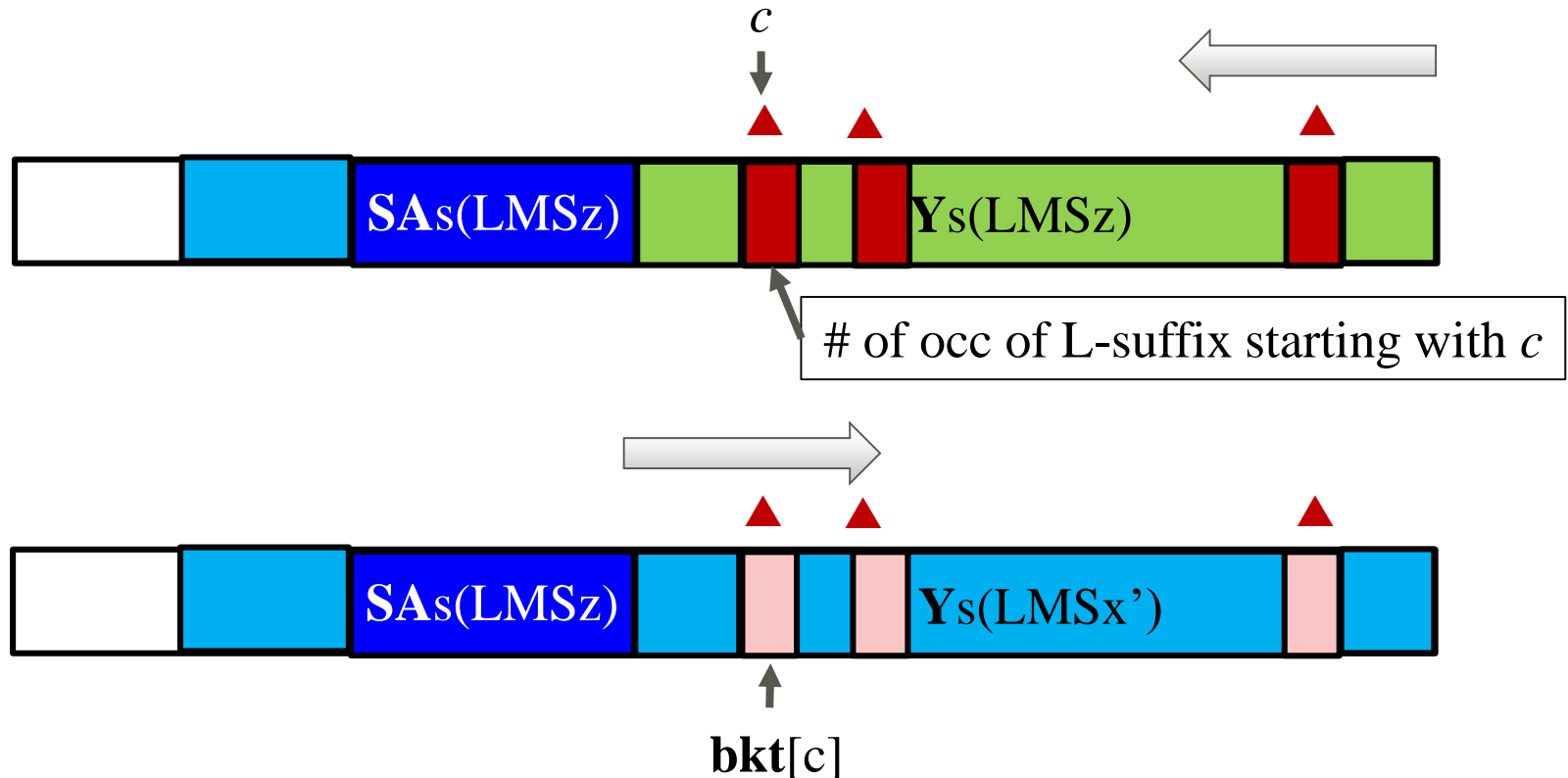
- If a suffix starting with c is in $s(LMSz')-s(LMSx)$, L-suffix starting with c exists
- By a right-to-left scan on T ,
 - We can determine $Y_{s(LMSz')}[T[i]]$
- We can move all suffixes of $s(LMSz')-s(LMSx)$ but they are not sorted
- We mark all all suffixes of $s(LMSz')-s(LMSx)$ and move them later by a sequential scan



How to Transit In-Place?

■ $\mathbf{bkt}_{s(\text{LMS}_x)}$:

- By a right-to-left scan on \mathbf{T} , count the number of L-suffixes starting with c , and store $\mathbf{Y}[c]$
- By a left-to-right scan on \mathbf{Y} , If the highest bit is 1, accumulate value, and decompress the value otherwise



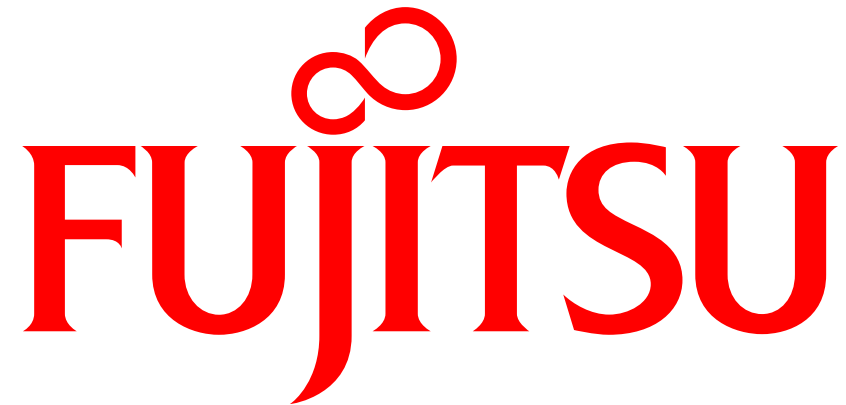
- If $|s(L)| > |s(S)|$,
 - Sort LML-suffixes
 - Sort S-suffixes
 - Sort L-suffixes

■ Summary

- We propose the first in-place suffix array construction algorithm
- We also propose a algorithm not to use stack space for recursive steps

■ Future Work

- More simpler in-place algorithm
- In-place linear time construction algorithm of Φ array
 - The most space efficient one is Goto and Bannai's algorithm which takes $\sigma \log N + O(\log N)$ bits of space
 - If it is possible, we can compute LZ77 in $N \log N$ bits + $O(\log N)$ bits of space.
the most space efficient one takes $N \log N + \sigma \log N + O(\log N)$ bits of space
[Goto and Bannai, 2014]



shaping tomorrow with you