



UNIVERSITÄT  
DES  
SAARLANDES



**ZBI**

ZENTRUM FÜR  
BIOINFORMATIK

# Linear Time Suffix Array Construction

Algorithms for Sequence Analysis

Sven Rahmann

Summer 2021

# Overview

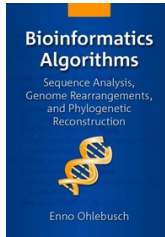
## Previous Lectures

- Ukkonen's algorithm: linear time **suffix tree** construction
- **Suffix links**
- Kasai's algorithm: linear time **LCP array** construction

## Today

- Direct linear time suffix array construction using **induced sorting**

# Recommended Literature

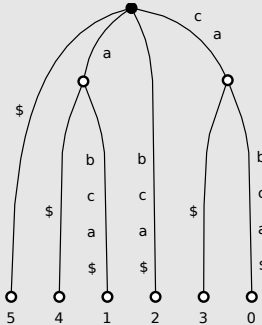


## Further Reading

- Shrestha et al. A bioinformatician's guide to the forefront of suffix array construction algorithms. Brief. Bioinformatics 2014 Mar;15(2):138-54
- G. Nong, S. Zhang and W. H. Chan. Linear Suffix Array Construction by Almost Pure Induced-Sorting. Proceedings of 19th Data Compression Conference (IEEE DCC), 2009.

# Suffix trees and suffix arrays

Suffix tree for the string  $T = cabca\$$ .

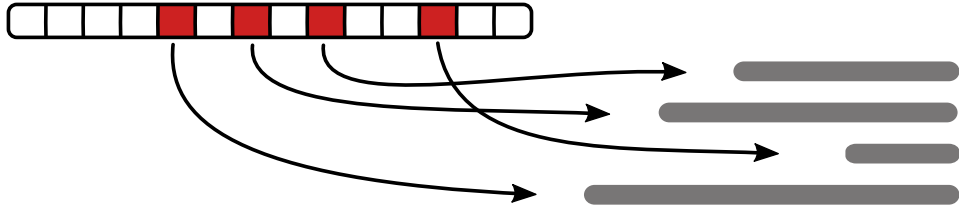


A suffix array of a string  $s\$$  with  $|s\$| = n$  is defined as the permutation  $pos$  of  $\{0, \dots, n-1\}$  that represents the lexicographic ordering of all suffixes of  $s\$$ .  
 $pos = [5, 4, 1, 2, 3, 0]$ .

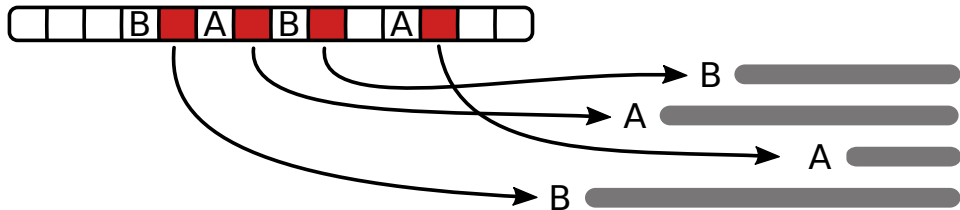
# Induced Sorting Idea



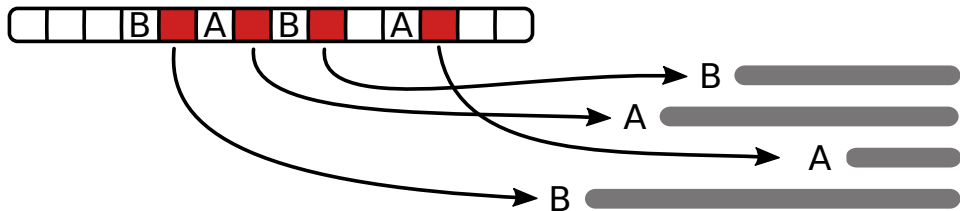
## Induced Sorting Idea



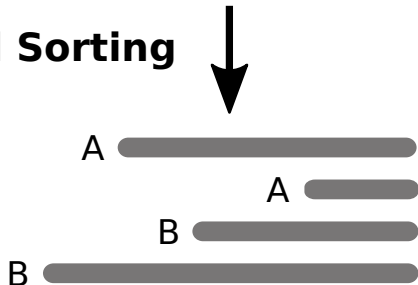
## Induced Sorting Idea



## Induced Sorting Idea



## Induced Sorting





# Definition of L-/S-positions

## Definition (L-position, S-position)

Let  $s\$$  be a string of length  $n$  with sentinel, such that  $s[n-1] = \$$ .

Let  $0 \leq p < n-1$  be a position in the text. We say,

- $p$  is an **L-position** (L means **larger**), if  $s[p \dots] > s[p+1 \dots]$ ,
- $p$  is an **S-position** (S means **smaller**), if  $s[p \dots] < s[p+1 \dots]$ ,
- The position of the sentinel  $n-1$  is defined as S-position.

(Note that no two suffixes can be identical.)

# Definition of L-/S-positions

## Definition (L-position, S-position)

Let  $s\$$  be a string of length  $n$  with sentinel, such that  $s[n-1] = \$$ .

Let  $0 \leq p < n-1$  be a position in the text. We say,

- $p$  is an **L-position** (L means **larger**), if  $s[p \dots] > s[p+1 \dots]$ ,
- $p$  is an **S-position** (S means **smaller**), if  $s[p \dots] < s[p+1 \dots]$ ,
- The position of the sentinel  $n-1$  is defined as S-position.

(Note that no two suffixes can be identical.)

	0.....1.....2.
Position $p$	0123456789012345678901
Sequence $s$	gccttaacattattacgccta\$
type	LSSLSSLSSLSSLSSLSSLLS

# Computing the L-/S-positions in the type array

The type information can be computed in linear time with a scan through the text from right to left:

```
1 def compute_types(T):
2     n = len(T)
3     typ = ['?'] * (n-1) + ['S']
4     for i in range(n-2, -1, -1):
5         typ[i] = 'L' if T[i] > T[i+1] else \
6                 'S' if T[i] < T[i+1] else typ[i+1]
7     return typ
```

In a real implementation, we use a bit vector (0/1) to represent the types.

# Definitions: LMS position / interval / substring

## Definition (LMS-interval, LMS-substring)

- S-positions located to the right of an L-position are called **LMS positions** (for **leftmost S** position).
- A pair of positions  $[i, j]$  is called **LMS interval** of  $s$ , if either
  - $i < j$  and both  $i$  and  $j$  are LMS-positions and there are no LMS-positions between  $i$  and  $j$ , or
  - $i = j = n - 1$ .
- Each LMS interval  $[i, j]$  is associated with its **LMS substring**  $s[i \dots j]$ .

## Observations

- Position  $n - 1$  with the sentinel is always an LMS-position.
- Whether an S-position is an LMS-position can be determined in constant time, looking up its type and the type to the left in the typearray.

## Example: type array, LMS substrings

	0	1	2
position p	0123456789012345678901		
sequence s	gccttaacattattacgccta\$		

## Example: type array, LMS substrings

	0	1	2
position p	0123456789012345678901		
sequence s	gccttaacattattacgccta\$		
type	LSSLSSLSSLSSLSSLSSL		

## Example: type array, LMS substrings

	0	1	2				
position p	0123456789012345678901						
sequence s	gccttaacattattacgccta\$						
type	LSSLSSLSSLSSLSSLSSL						
LMS?	*	*	*	*	*	*	*
LMS-substr	cctta	atta	acgc	\$			
		aaca	atta	ccta\$			

# Overview of Induced Sorting

## Notation

- $s$  is the input sequence,
- $pos$  is the desired output suffix array of  $s$ .

## Induced sorting

- Scan  $s$  to compute the type array
- Scan type to find all LMS positions in  $s$
- Phase I - Sort suffixes at LMS positions (complex; recursive)
- Phase II - Sort all remaining suffixes of  $s$  (easy)
- Output  $pos$



## Code: Overview

```
1 def sais_main(T, alphabet_size):
2     # T: text (bytes, numpy array, not str!), T[n-1]=0
3     # alphabet_size, 1 <= T[i] < alphabet_size for all i < n-1
4
5     pos = np.empty(len(T), dtype=np.int64)
6     # B[a]: total number of characters in T that are <= a
7     B = count_cumulative_characters(T, alphabet_size)
8     types = compute_types(T)
9     lms_positions = find_lms_positions(types)
10    # Phase 1 sorts lms_positions lexicographically in-place,
11    # may recurse into sais_main() with a reduced text.
12    phase1(T, B, types, lms_positions, pos)
13    # Phase 2 sorts all suffixes from correctly sorted LMS.
14    phase2(T, B, types, lms_positions, pos)
15    return pos
```

## Code: Initialization, buckets and types

```
1 def count_cumulative_characters(T, alphabet_size):
2     # B[a]: total number of characters in T that are <= a
3     B = np.zeros(alphabet_size, dtype=np.uint64)
4     for a in T:
5         B[a] += 1
6     for a in range(1, alphabet_size):
7         B[a] += B[a-1]
8     return B
```

```
1 def compute_types(T):
2     # Compute position types (SMALLER=0, LARGER=1) for T
3     n = len(T)
4     types = np.zeros(n, dtype=np.uint8) # types[n-1] = SMALLER
5     for i in range(n-2, -1, -1):
6         types[i] = LARGER if T[i] > T[i+1] else \
7                     SMALLER if T[i] < T[i+1] else types[i+1]
8     return types
```

## Code: Initialization, LMS positions

```
1 def find_lms_positions(types):
2     n = len(types)
3     # count the number of LMS positions first
4     m = 0
5     for p in range(1, n):
6         m += (types[p] == SMALLER and types[p-1] == LARGER)
7     # allocate array of just the correct size m
8     lms_positions = np.empty(m, dtype=np.int64)
9     # now fill the array with the actual LMS positions
10    m = 0
11    for p in range(1, n):
12        if types[p] == SMALLER and types[p-1] == LARGER:
13            lms_positions[m] = p
14            m += 1
15    return lms_positions
```

## Code: Overview again

```
1 def sais_main(T, alphabet_size):
2     # T: text (bytes, numpy array, not str!), T[n-1]=0
3     # alphabet_size, 1 <= T[i] < alphabet_size for all i < n-1
4
5     pos = np.empty(len(T), dtype=np.int64)
6     # B[a]: total number of characters in T that are <= a
7     B = count_cumulative_characters(T, alphabet_size)
8     types = compute_types(T)
9     lms_positions = find_lms_positions(types)
10    # Phase 1 sorts lms_positions lexicographically in-place,
11    # may recurse into sais_main() with a reduced text.
12    phase1(T, B, types, lms_positions, pos)
13    # Phase 2 sorts all suffixes from correctly sorted LMS.
14    phase2(T, B, types, lms_positions, pos)
15    return pos
```

# Phase II

# Sorting the non-LMS suffixes

Let's start with Phase II (Phase I uses elements of Phase II):

## Definition (Bucket)

A maximal interval of the suffix array  $pos$ , in which the referenced suffixes start with the same character, is called a **bucket**.

There are as many buckets as characters in the used alphabet, plus the one for the sentinel character.

# Sorting the non-LMS suffixes

## Lemma

*Within each bucket of the suffix array, the L-positions appear before the S-positions.*

## Proof

Let  $p$  be an S-position, and let  $q$  be an L-position, let  $s[p] = s[q] = b \in \Sigma$ , so both  $p$  and  $q$  are in the  $b$ -bucket. Then the suffix  $p + 1$  is **larger** than suffix  $p$ , and suffix  $q + 1$  is **smaller** than suffix  $q$ . Because  $s[p] = s[q]$ , the order of  $p$  vs.  $q$  is determined by  $p + 1$  vs.  $q + 1$ , but  $q + 1$  comes before  $p + 1$  in the lexicographic order.

## Illustration

Let  $a < b < c$ ; suffix  $q$  is  $b^+a$ , whereas  $p$  is  $b^+c$ :

q		p
bbb...a	<	bbb...c
L		S

# Sorting the non-LMS suffixes

## Lemma

*Within each bucket of the suffix array, the L-positions appear before the S-positions.*

Bucket	\$		a		c		t	
	L	S	L	S	L	S	L	S



# Sorting the non-LMS suffixes

## Idea

- Use the already sorted **LMS-positions** (a subset of the S-positions) to **sort the L-positions** correctly, and then
- use the sorted **L-positions** to sort all **S-positions**.

This is why the algorithm is called **induced sorting**:

The order of one type of suffixes completely induces the ordering of the others.

# Preparing the Suffix Array

## Step (1)

- Initialize pos with **unknown** at each position
- Mark the beginning and end of each bucket by pointers
- Write the **sorted** LMS-positions (phase I) at the end of their respective buckets.

# Preparing the Suffix Array

## Step (1)

- Initialize pos with **unknown** at each position
- Mark the beginning and end of each bucket by pointers
- Write the **sorted** LMS-positions (phase I) at the end of their respective buckets.

	0	1	2
position p	0123456789012345678901		
sequence s	gccttaacattattacgccta\$		
type	LSSLLSSL	LSSLLSSL	LSSLLSSL
LMS?	*	*	*

rank r	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21
bucket	\$	a	a	a	a	a	a	c	c	c	c	c	c	g	g	t	t	t	t	t	t	t
pos/	21	.	.	5	14	11	8	.	.	.	.	17	1	.	.		.	.	.	.	.	

# Sorting the L-positions (Induced Sorting)

## Step (2)

- Iterate through  $\text{pos}$  from **left to right** with index  $r$ .
- If  $\text{pos}[r]$  is unknown, skip index  $r$ .
- Otherwise, look at  $\text{pos}[r] - 1$ :
  - 1 If  $\text{pos}[r] - 1$  is an L-position, enter it at the first free position in its bucket.
  - 2 If  $\text{pos}[r] - 1$  is an S-position, skip index  $r$ .

## Result

All L-positions are entered in the suffix array in correct order.

# Example: Sorting the L-positions

	0	1	2
position p	0123456789012345678901		
sequence s	gccttaacattattacgccta\$		
type	LSSLLSSLSLLSLLSSLLS		
LMS?	*	*	*

rank r	0  1	2	3	4	5	6  7	8	9	10	11	12 13	14 15	16	17	18	19	20	21
bucket	\$  a	a	a	a	a	a  c	c	c	c	c	c  g	g  t	t	t	t	t	t	t
pos	21  .	.	5	14	11	8  .	.	.	.	17	1  .	.	.	.	.	.	.	.
	^S  vL																	
pos	21 20	.	5	14	11	8  .	.	.	.	17	1  .	.	.	.	.	.	.	.
	^L												vL					
pos	21 20	.	5	14	11	8  .	.	.	.	17	1  .	.	19	.	.	.	.	.
			^S										vL					
pos	21 20	.	5	14	11	8  .	.	.	.	17	1  .	.	19	4	.	.	.	.
				^S									vL					
pos	21 20	.	5	14	11	8  .	.	.	.	17	1  .	.	19	4	13	.	.	.
					^S										vL			... ..
pos	21 20	.	5	14	11	8  7	.	.	.	17	1 16	0 19	4	13	10	3	12	9

# Sorting the S-positions (Induced Sorting)

## Step (3)

- 1 Remove all the S-positions from  $\text{pos}$ , except \$.
- 2 Iterate through  $\text{pos}$  from **right to left** with index  $r$ .
- 3 If  $\text{pos}[r]$  is unknown, skip index  $r$ .
- 4 Otherwise, look at  $\text{pos}[r] - 1$ :
  - If  $\text{pos}[r] - 1$  is an S-position, enter it at the rightmost free position in its bucket.
  - If  $\text{pos}[r] - 1$  is an L-position, skip index  $r$ .

## Result

All S-positions are entered in the suffix array in correct order.

# Example: Sorting the S-positions (Induced Sorting)

	0	1	2
position p	0123456789012345678901		
sequence s	gccttaacattattacgccta\$		
type	LSSLLSSLSLLSLLSSLLS		
LMS?	* * * * *		

rank r	0  1	2	3	4	5	6  7	8	9	10	11	12 13	14 15	16	17	18	19	20	21
bucket	\$  a	a	a	a	a	a  c	c	c	c	c	c  g	g  t	t	t	t	t	t	t
pos(2)	21 20	.	5	14	11	8  7	.	.	.	17	1 16	0 19	4	13	10	3	12	9

pos	21 20	.	.	.	.	8  7	.	.	.	.	1 16	0 19	4	13	10	3	12	9
pos	21 20	.	.	.	11	8  7	.	.	.	.	1 16	0 19	4	13	10	3	12	9
pos	21 20	.	.	.	11	8  7	.	.	.	.	2 16	0 19	4	13	10	3	12	9
pos	21 20	.	.	.	11	8  7	.	.	.	18	2 16	0 19	4	13	10	3	12	9
pos	21 20	.	.	.	11	8  7	.	.	15	18	2 16	0 19	4	13	10	3	12	9
pos	21 20	.	.	.	11	8  7	.	1	15	18	2 16	0 19	4	13	10	3	12	9
pos	21 20	.	.	.	11	8  7	17	1	15	18	2 16	0 19	4	13	10	3	12	9
pos	21 20	.	.	14	11	8  7	17	1	15	18	2 16	0 19	4	13	10	3	12	9
pos	21 20	.	6	14	11	8  7	17	1	15	18	2 16	0 19	4	13	10	3	12	9
pos	21 20	5	6	14	11	8  7	17	1	15	18	2 16	0 19	4	13	10	3	12	9
pos	21 20	5	6	14	11	8  7	17	1	15	18	2 16	0 19	4	13	10	3	12	9

# Summary and Analysis of Phase II

## Phase II:

- 1 Enter sorted LMS suffixes into pos, set bucket pointers
- 2 Sort L-suffixes based on sorted LMS-suffixes (induced sorting)
- 3 Sort S-suffixes based on sorted L-suffixes (induced sorting)



# Summary and Analysis of Phase II

## Phase II:

- 1 Enter sorted LMS suffixes into pos, set bucket pointers
- 2 Sort L-suffixes based on sorted LMS-suffixes (induced sorting)
- 3 Sort S-suffixes based on sorted L-suffixes (induced sorting)

## Running Time Analysis

- Step (1) can be done in linear time.
- Step (2) and (3) each do a linear scan through the suffix array in linear time.
- $\Rightarrow$  Phase II takes linear time.

# Summary and Analysis of Phase II

## Phase II:

- 1 Enter sorted LMS suffixes into pos, set bucket pointers
- 2 Sort L-suffixes based on sorted LMS-suffixes (induced sorting)
- 3 Sort S-suffixes based on sorted L-suffixes (induced sorting)

## Running Time Analysis

- Step (1) can be done in linear time.
- Step (2) and (3) each do a linear scan through the suffix array in linear time.
- $\Rightarrow$  Phase II takes linear time.

## Correctness?

# Correct Sorting of L-Positions

## Lemma: Correctness of Step (2)

Assuming correctly ordered LMS-positions in each bucket, then after Step (2), **all** L-positions can be found at their **correct** positions.

# Correct Sorting of L-Positions

## Lemma: Correctness of Step (2)

Assuming correctly ordered LMS-positions in each bucket, then after Step (2), **all** L-positions can be found at their **correct** positions.

## Proof idea

- If  $p$  is a text position with rank  $r$  in  $\text{pos}$  and  $p - 1$  is a L-position, then  $p - 1$  has a rank  $r'$  with  $r' > r$  by definition of an L-position.

# Correct Sorting of L-Positions

## Lemma: Correctness of Step (2)

Assuming correctly ordered LMS-positions in each bucket, then after Step (2), **all** L-positions can be found at their **correct** positions.

## Proof idea

- If  $p$  is a text position with rank  $r$  in  $\text{pos}$  and  $p - 1$  is a L-position, then  $p - 1$  has a rank  $r'$  with  $r' > r$  by definition of an L-position.
- This assures that each L-position  $p - 1$  will
  - 1 be induced by an LMS- or L-position  $p$
  - 2 be induced by a position further to the left

# Correct Sorting of L-Positions

## Lemma: Correctness of Step (2)

Assuming correctly ordered LMS-positions in each bucket, then after Step (2), **all** L-positions can be found at their **correct** positions.

## Proof idea

- If  $p$  is a text position with rank  $r$  in  $\text{pos}$  and  $p - 1$  is a L-position, then  $p - 1$  has a rank  $r'$  with  $r' > r$  by definition of an L-position.
- This assures that each L-position  $p - 1$  will
  - 1 be induced by an LMS- or L-position  $p$
  - 2 be induced by a position further to the left
- Complete proof by induction:  
Show that the first  $k$  LMS- and L-positions all appear in the correct order.

# Correct Sorting of S-Positions

## Lemma: Correctness of Step (3)

Assuming correctly ordered L-positions in each bucket, then after step (3), **all** positions can be found at their **correct** positions.

# Correct Sorting of S-Positions

## Lemma: Correctness of Step (3)

Assuming correctly ordered L-positions in each bucket, then after step (3), **all** positions can be found at their **correct** positions.

## Proof idea

- Let  $p$  be a text position with rank  $r$  in pos and  $p - 1$  is a S-position, then  $p - 1$  has a rank  $r'$  with  $r' < r$  (by definition of an S-position).



# Correct Sorting of S-Positions

## Lemma: Correctness of Step (3)

Assuming correctly ordered L-positions in each bucket, then after step (3), **all** positions can be found at their **correct** positions.

## Proof idea

- Let  $p$  be a text position with rank  $r$  in pos and  $p - 1$  is a S-position, then  $p - 1$  has a rank  $r'$  with  $r' < r$  (by definition of an S-position).
- This assures that each S-position  $p - 1$  will be induced by a position  $p$  further to the right.

# Correct Sorting of S-Positions

## Lemma: Correctness of Step (3)

Assuming correctly ordered L-positions in each bucket, then after step (3), **all** positions can be found at their **correct** positions.

## Proof idea

- Let  $p$  be a text position with rank  $r$  in pos and  $p - 1$  is a S-position, then  $p - 1$  has a rank  $r'$  with  $r' < r$  (by definition of an S-position).
- This assures that each S-position  $p - 1$  will be induced by a position  $p$  further to the right.
- Complete proof by induction (in  $k$ ):  
Show that the last  $k$  positions all appear in the correct order.

## Code: Phase II

```
1 def phase2(T, B0, types, lms, pos):
2     # T: Text, B0: cumulative bucket sizes, types: type array
3     # lms: sorted or unsorted LMS positions
4     # pos: suffix array (output)
5
6     # 0. Initialize pos by inserting LMS positions,
7     B = B0.copy() # working copy of C, to be modified
8     initialize_pos_from_lms(T, B, lms, pos)
9     # 1. Do a left-to-right induction scan for L-positions,
10    B[:] = B0[:] # re-set B to a clean working copy of C
11    induce_L_positions(T, B, types, pos)
12    # 2. Do a right-to-left induction scan for S-positions.
13    B[:] = B0[:] # re-set B to a clean working copy of C
14    induce_S_positions(T, B, types, pos)
15    # Result: pos has been modified as described.
```

## Code: Phase II, Initialization

```
1 def initialize_pos_from_lms(T, B, lms, pos):  
2     pos[:] = -1 # set everything to "unknown"  
3     # Insert LMS positions at right end of their buckets,  
4     # right-to-left, so we know where to start in each bucket.  
5     for p in lms[::-1]:  
6         a = T[p] # character determines the bucket  
7         B[a] -= 1  
8         pos[B[a]] = p
```

## Code: Phase II, L-positions

```
1 def induce_L_positions(T, B, types, pos):
2     # Left-to-right scan: Induce L-positions from LMS-positions
3     n = len(T)
4     for r in range(n):
5         p = pos[r]
6         if p <= 0: continue # unknown or 0 -> skip
7         if types[p-1] == SMALLER: continue # skip S positions
8         a = T[p-1] # determine bucket
9         pos[B[a-1]] = p-1
10        B[a-1] += 1
```

## Code: Phase II, S-positions

```
1 def induce_S_positions(T, B, types, pos):
2     # Right-to-left scan: Induce S-positions from L-positions
3     n = len(T)
4     for r in range(n-1, -1, -1):
5         p = pos[r]
6         if p == 0: continue # skip position 0 (no p-1)
7         if types[p-1] == LARGER: continue # skip L positions
8         a = T[p-1] # determine bucket
9         B[a] -= 1
10        pos[B[a]] = p-1
```

# Phase I

# Idea for Phase 1

## Goal (hard)

Sort the LMS suffixes (i.e., suffixes starting at LMS positions)



# Idea for Phase 1

## Goal (hard)

Sort the LMS suffixes (i.e., suffixes starting at LMS positions)

## Plan

- Only sort the LMS substrings (up to next LMS position): shorter total length ( $O(n)$  instead of  $O(n^2)$ ).
- Expand alphabet and reduce text length (LMS substring  $\mapsto$  character), keeping lexicographic order of LMS substrings (“**lexicographic naming**”).
- If all LMS substrings are distinct, we have also sorted the LMS suffixes, done!
- If there are equal LMS substrings, compute suffix array of reduced text (recursively with SAIS), use that to infer correct order of LMS suffixes.

# Example: Alphabet Expansion and Text Reduction

	0	1	2
position p	0123456789012345678901		
text T	gccttaacattattacgccta\$		
type	LSSLSSLSSLSSLSSLSSL		
LMS?	* * * * *		
LMS-substr	cctta	atta	acgc \$
	aaca	atta	ccta\$

	0	1	2	3	4	5	6
p'							
red. text R	E	A	C	C	B	D	\$

	0	1	2	3	4	5	6	
r'								
pos' [r']	6	1	4	3	2	5	0	reduced suffix array
RT[pos' [r']]	21	5	14	11	8	17	1	sorted LMS-positions

# Overview with Recursion

- 1) **Phase 1:** Identify and sort LMS substrings
- 2) Reduce text by lexicographic naming

A=aaca, B=acgc, ...

E A C C B D \$

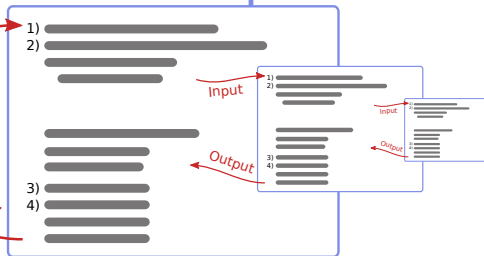
Input

Suffix array of  
reduced text

[6 1 4 3 2 5 0]

Output

- 3) Translate into sorted  
LMS suffixes
- 4) **Phase 2:** Use sorted LMS suffixes  
to induce order of non-LMS suffixes



# Achieving Phase I in Linear Time

## Questions

- 1 How to sort the LMS substrings in linear time?
- 2 How to compare and name LMS substrings in linear time?
- 3 How to obtain order of LMS suffixes after recursive call ?

# Achieving Phase I in Linear Time

## Questions

- 1 How to sort the LMS substrings in linear time?
- 2 How to compare and name LMS substrings in linear time?
- 3 How to obtain order of LMS suffixes after recursive call ?

## Sorting LMS Substrings

Surprisingly, it can be done by another run of Phase II:

- Enter **unsorted** LMS-positions into correct buckets of pos
- Induce order of L-positions based on **unsorted** LMS-positions
- Induce order of S-positions based on sorted S-positions

**Result:** Suffixes at LMS positions correctly sorted **up to next LMS position:**

... SSSLLLLLS ...  
... \* \* ...

# Achieving Phase I in Linear Time

## Text Reduction and Lexicographic Naming

- 1 Sort LMS substrings (phase 2) into pos (previous slide)
- 2 Extract partially sorted LMS positions from pos
- 3 Compare LMS substrings in lexicographic order, \$ first, assign new “name” (number) if different from previous string.
- 4 In parallel, build new reduced text  $R$  from names at LMS positions, build map  $RT$  from  $R$ -positions to  $T$ -LMS-positions.
- 5 If all LMS substrings are unique, we already have sorted LMS suffixes. Otherwise recurse on  $R$  (next slide) to obtain  $pos'$ .
- 6 Total time without recursion:  $O(n)$ .

# Achieving Phase I in Linear Time

## Recursion

**Situation:** We have

- paritally sorted LMS suffixes  $lms$ ,
- reduced text  $R$ ,
- map  $RT$  from  $R$ -positions to  $T$ -LMS-positions.

**Left to do:**

- 1 Recursively compute  $pos'$  of  $R$  by calling  $SAIS(R)$ .
- 2 Overwrite  $lms$  by correct order of  $T$  is  $RT[pos'[0]], RT[pos'[1]], \dots$

## Code: Phase I, Overview

```
1 def phase1(T, B, types, lms_positions, pos):
2     # T: text; B: cumulative character counts
3     # lms_positions: LMS positions in ANY ORDER
4     # pos: uninitialized, used to sort LMS positions
5     alphabet_size = len(B)
6     phase2(T, B, types, lms_positions, pos)
7     # Compute reduced text from LMS substrings
8     (R, reduced_alphabet_size, position_map) \
9     = reduce_text(T, alphabet_size, types, pos, lms_positions)
10    # If there are equal LMS substrings, recurse on reduced text
11    if len(R) != reduced_alphabet_size:
12        reduced_pos = sais_main(R, reduced_alphabet_size)
13        # Re-map reduced_pos to original text positions;
14        # these are the lms_positions in lexicographic order,
15        for i, redp in enumerate(reduced_pos):
16            lms_positions[i] = position_map[redp]
```



## Code: Phase I, Text Reduction (Lexicographic Naming)

```
1 def reduce_text(T, alphabet_size, types, pos, lms_positions):
2     n, m = len(pos), len(lms_positions)
3     names = np.full(n, -1, dtype=np.int64) # the names
4     last_lms = n-1; names[last_lms] = 0 # sentinel at n-1
5     reduced_alphabet_size = 1; j = 0
6     # go through the suffixes lexicographically, w/o sentinel
7     for r in range(1, n):
8         p = pos[r] # if not LMS, skip it:
9         if p==0 or types[p]!=SMALLER or types[p-1]!=LARGER:
10             continue
11         lms_positions[j]=p; j+=1 # write sorted LMS positions
12         if lms_substrings_unequal(T, types, last_lms, p):
13             reduced_alphabet_size += 1
14         names[p] = reduced_alphabet_size - 1
15         last_lms = p
```

## Code: Phase I, Comparison of LMS Substrings

```
1 def lms_substrings_unequal(T, types, p1, p2):
2     """Return True iff LMS substrings at p1, p2 in T differ"""
3     is_lms_p1 = is_lms_p2 = False
4     while True:
5         if T[p1] != T[p2]: return True # unequal
6         if types[p1] != types[p2]: return True # unequal
7         if is_lms_p1 and is_lms_p2: return False # equal
8         p1 +=1; p2 += 1 # look at next positions
9         # check if both or only one LMS substring ends now
10        is_lms_p1 = types[p1]==SMALLER and types[p1-1]==LARGER
11        is_lms_p2 = types[p2]==SMALLER and types[p2-1]==LARGER
12        if is_lms_p1 and is_lms_p2: continue # final test
13        if is_lms_p1 or is_lms_p2: return True # unequal
```

# Running Time Analysis

## Observations about the recursion

- The alphabet size can grow, but is bounded by  $n$  (e.g. a, c, g, t expands to A–E).
- After each reduction step for a sequence of length  $n$  (including the sentinel), the new sequence has at most length  $\lfloor n/2 \rfloor$  (again including the sentinel).

# Running Time Analysis

## Observations about the recursion

- The alphabet size can grow, but is bounded by  $n$  (e.g. a, c, g, t expands to A–E).
- After each reduction step for a sequence of length  $n$  (including the sentinel), the new sequence has at most length  $\lfloor n/2 \rfloor$  (again including the sentinel).

Find a bound on the running time  $T(n)$  for these three parts:

- 1 Phase I without recursion:  $\leq c_1 n$
- 2 Recursive call:  $\leq T(n/2)$
- 3 Phase II:  $\leq c_2 n$

# Running Time Analysis

## Observations about the recursion

- The alphabet size can grow, but is bounded by  $n$  (e.g. a, c, g, t expands to A–E).
- After each reduction step for a sequence of length  $n$  (including the sentinel), the new sequence has at most length  $\lfloor n/2 \rfloor$  (again including the sentinel).

Find a bound on the running time  $T(n)$  for these three parts:

- 1 Phase I without recursion:  $\leq c_1 n$
- 2 Recursive call:  $\leq T(n/2)$
- 3 Phase II:  $\leq c_2 n$

## Claim

$T(n) = \mathcal{O}(n)$ , i.e., SAIS takes linear time in  $n = |T|$ .

# Running Time Analysis (Proof)

## Proof of Claim $T(n) = \mathcal{O}(n)$

- 1 Phase I without recursion:  $\leq c_1 n$
- 2 Recursive call:  $\leq T(n/2)$
- 3 Phase II:  $\leq c_2 n$

Let  $C := c_1 + c_2$ . Then  $T(1) = \mathcal{O}(1)$ , and thus

$$\begin{aligned} T(n) &\leq c_1 n + T(n/2) + c_2 n \\ &= C n + T(n/2) \\ &= C n + C n/2 + T(n/4) \\ &\leq C n(1 + 1/2 + 1/4 + \dots) + T(1) \\ &= 2C n + \mathcal{O}(1) = \mathcal{O}(n). \end{aligned}$$

q.e.d.

# Summary

## Linear suffix array construction by **induced sorting** (SAIS)

- 1 Sorted LMS-suffixes can be used to induce sorting of L-suffixes.
- 2 Sorted L-suffixes can be used to induce sorting of S-suffixes.
- 3 Sort LMS-suffixes by sorting LMS-substrings first  
(how? induced sorting on **unsorted** LMS-positions)
- 4 Reduce text by lexicographic naming of LMS-substrings
- 5 If equal LMS-substrings exist, recurse on reduced text
- 6 LMS-order of original text is obtained from suffix array of reduced text

# Possible exam questions

- Explain the principle of induced sorting.
- Why are L-positions on the left and S-positions on the right of each bucket?
- What is the goal of the text reduction step?
- Conduct the first iteration of induced sorting for a small example string.
- Explain why the induced sorting algorithm has linear running time.