

Longest Common Extension Queries via String Synchronizing Sets

Tomasz Kociumaka

Randomized Algorithms and Probabilistic Analysis of Algorithms

- [1] Tomasz Kociumaka. Efficient Data Structures for Internal Queries in Texts. PhD thesis, University of Warsaw, 2018.
- [2] Dominik Kempa and Tomasz Kociumaka. String synchronizing sets: Sublinear-time BWT construction and optimal LCE data structure. STOC 2019.
- [3] Patrick Dinklage, Johannes Fischer, Alexander Herlez, Tomasz Kociumaka, and Florian Kurpicz. Practical Performance of Space Efficient Data Structures for Longest Common Extensions. ESA 2020.

February 8, 2023

Longest Common Extension Queries

Definition

The *Longest Common Extension* $LCE(i, j)$ of positions i, j of a string T (of length n) is the length of the longest common prefix of the suffixes $T[i..n]$ and $T[j..n]$ of T .

Longest Common Extension Queries

Definition

The *Longest Common Extension* $LCE(i, j)$ of positions i, j of a string T (of length n) is the length of the longest common prefix of the suffixes $T[i..n]$ and $T[j..n]$ of T .

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

Longest Common Extension Queries

Definition

The *Longest Common Extension* $LCE(i, j)$ of positions i, j of a string T (of length n) is the length of the longest common prefix of the suffixes $T[i..n]$ and $T[j..n]$ of T .

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

$LCE(16, 7) =$

Longest Common Extension Queries

Definition

The *Longest Common Extension* $LCE(i, j)$ of positions i, j of a string T (of length n) is the length of the longest common prefix of the suffixes $T[i..n]$ and $T[j..n]$ of T .

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

The diagram shows a binary string of length 20: 1 1 0 1 0 0 1 0 1 0 1 0 1 0 0 1 0 1 0 0. The characters are indexed from 1 to 20. Two suffixes are highlighted with green boxes and arrows: one starting at index 7 (1 0 1 0) and another starting at index 16 (1 0 1 0). The character at index 11 (1) is highlighted in red, and the character at index 20 (0) is also highlighted in red.

$LCE(16, 7) =$

Longest Common Extension Queries

Definition

The *Longest Common Extension* $LCE(i, j)$ of positions i, j of a string T (of length n) is the length of the longest common prefix of the suffixes $T[i..n]$ and $T[j..n]$ of T .

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

$$LCE(16, 7) = 4$$

Longest Common Extension Queries

Definition

The *Longest Common Extension* $LCE(i, j)$ of positions i, j of a string T (of length n) is the length of the longest common prefix of the suffixes $T[i..n]$ and $T[j..n]$ of T .

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

$$LCE(16, 7) = 4$$

Used as a *subroutine* in many algorithms and data structures such as:

- approximate pattern matching (the *kangaroo method*),
- discovery of repetitions in strings,
- construction of text indexing data structures.

Classic Data Structure for LCE Queries

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

Components:

Classic Data Structure for LCE Queries

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

Components:

- 1 Lexicographic rank of each suffix,

<i>i</i>	
1	0
2	00
3	0010100
4	0010101010010100
5	0100
6	010010100
7	010010101010010100
8	010100
9	01010010100
10	0101010010100
11	010101010010100
12	100
13	10010100
14	10010101010010100
15	10100
16	1010010100
17	1010010101010010100
18	101010010100
19	10101010010100
20	11010010101010010100

Classic Data Structure for LCE Queries

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

Components:

- 1 Lexicographic rank of each suffix,

<i>i</i>	
1	0
2	00
3	0010100
4	0010101010010100
5	0100
6	010010100
7	010010101010010100
8	010100
9	01010010100
10	0101010010100
11	010101010010100
12	100
13	10010100
14	10010101010010100
15	10100
16	1010010100
17	1010010101010010100
18	101010010100
19	10101010010100
20	11010010101010010100

Classic Data Structure for LCE Queries

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

Components:

- 1 Lexicographic rank of each suffix,
- 2 LCP lengths for subsequent suffixes.

LCP[i]	i	
-	1	0
1	2	00
2	3	0010100
6	4	0010101010010100
1	5	0100
4	6	010010100
8	7	010010101010010100
3	8	010100
6	9	01010010100
5	10	0101010010100
7	11	010101010010100
0	12	100
3	13	10010100
7	14	10010101010010100
2	15	10100
5	16	1010010100
9	17	1010010101010010100
4	18	101010010100
6	19	10101010010100
1	20	11010010101010010100

Classic Data Structure for LCE Queries

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

Components:

- 1 Lexicographic rank of each suffix,
- 2 LCP lengths for subsequent suffixes.

LCP[i]	i	
-	1	0
1	2	00
2	3	0010100
6	4	0010101010010100
1	5	0100
4	6	010010100
8	7	010010101010010100
3	8	010100
6	9	01010010100
5	10	0101010010100
7	11	010101010010100
0	12	100
3	13	10010100
7	14	10010101010010100
2	15	10100
5	16	1010010100
9	17	1010010101010010100
4	18	101010010100
6	19	10101010010100
1	20	11010010101010010100

Classic Data Structure for LCE Queries

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

Components:

- 1 Lexicographic rank of each suffix,
- 2 LCP lengths for subsequent suffixes.

Queries: $LCE(i, j) = \min LCP(\text{rank}[i].. \text{rank}[j])$

LCP[i]	i	
-	1	0
1	2	00
2	3	0010100
6	4	0010101010010100
1	5	0100
4	6	010010100
8	7	010010101010010100
3	8	010100
6	9	01010010100
5	10	0101010010100
7	11	010101010010100
0	12	100
3	13	10010100
7	14	10010101010010100
2	15	10100
5	16	1010010100
9	17	1010010101010010100
4	18	101010010100
6	19	10101010010100
1	20	11010010101010010100

Classic Data Structure for LCE Queries

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

Components:

- 1 Lexicographic rank of each suffix,
- 2 LCP lengths for subsequent suffixes.

Queries: $LCE(i, j) = \min LCP(\text{rank}[i] \dots \text{rank}[j])$

$LCE(16, 7) =$

LCP[i]	i	
-	1	0
1	2	00
2	3	0010100
6	4	0010101010010100
1	5	0100
4	6	010010100
8	7	010010101010010100
3	8	010100
6	9	01010010100
5	10	0101010010100
7	11	010101010010100
0	12	100
3	13	10010100
7	14	10010101010010100
2	15	10100
5	16	1010010100
9	17	1010010101010010100
4	18	101010010100
6	19	10101010010100
1	20	11010010101010010100

Classic Data Structure for LCE Queries

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

Components:

- 1 Lexicographic rank of each suffix,
- 2 LCP lengths for subsequent suffixes.

Queries: $LCE(i, j) = \min LCP(\text{rank}[i] .. \text{rank}[j])$

$LCE(16, 7) =$

LCP[i]	i	
-	1	0
1	2	00
2	3	0010100
6	4	0010101010010100
1	5	0100
4	6	010010100
8	7	010010101010010100
3	8	010100
6	9	01010010100
5	10	0101010010100
7	11	010101010010100
0	12	100
3	13	10010100
7	14	10010101010010100
2	15	10100
5	16	1010010100
9	17	1010010101010010100
4	18	101010010100
6	19	10101010010100
1	20	11010010101010010100

Classic Data Structure for LCE Queries

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

Components:

- 1 Lexicographic rank of each suffix,
- 2 LCP lengths for subsequent suffixes.

Queries: $LCE(i, j) = \min LCP(\text{rank}[i] \dots \text{rank}[j])$

$LCE(16, 7) =$

LCP[i]	i	
-	1	0
1	2	00
2	3	0010100
6	4	0010101010010100
1	5	0100
4	6	010010100
8	7	010010101010010100
3	8	010100
6	9	01010010100
5	10	0101010010100
7	11	010101010010100
0	12	100
3	13	10010100
7	14	10010101010010100
2	15	10100
5	16	1010010100
9	17	1010010101010010100
4	18	101010010100
6	19	10101010010100
1	20	11010010101010010100

Classic Data Structure for LCE Queries

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

Components:

- 1 Lexicographic rank of each suffix,
- 2 LCP lengths for subsequent suffixes.

Queries: $LCE(i, j) = \min LCP(\text{rank}[i] \dots \text{rank}[j])$

$$LCE(16, 7) = 4$$

LCP[i]	i	
-	1	0
1	2	00
2	3	0010100
6	4	0010101010010100
1	5	0100
4	6	010010100
8	7	010010101010010100
3	8	010100
6	9	01010010100
5	10	0101010010100
7	11	010101010010100
0	12	100
3	13	10010100
7	14	10010101010010100
2	15	10100
5	16	1010010100
9	17	1010010101010010100
4	18	101010010100
6	19	10101010010100
1	20	11010010101010010100

Classic Data Structure for LCE Queries

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

Components:

- 1 Lexicographic rank of each suffix,
- 2 LCP lengths for subsequent suffixes.

Queries: $LCE(i, j) = \min LCP(\text{rank}[i].. \text{rank}[j])$

$$LCE(16, 7) = 4$$

Efficiency:

- query time: $\mathcal{O}(1)$;
- data structure size: $\mathcal{O}(n)$;
- construction time: $\mathcal{O}(n)$.

LCP[i]	i	
-	1	0
1	2	00
2	3	0010100
6	4	0010101010010100
1	5	0100
4	6	010010100
8	7	010010101010010100
3	8	010100
6	9	01010010100
5	10	0101010010100
7	11	010101010010100
0	12	100
3	13	10010100
7	14	10010101010010100
2	15	10100
5	16	1010010100
9	17	1010010101010010100
4	18	101010010100
6	19	10101010010100
1	20	11010010101010010100

Is $\mathcal{O}(n)$ space optimal?

Is $\mathcal{O}(n)$ space optimal?

- The string T can be encoded in $\mathcal{O}(n \log \sigma)$ bits, where σ is the alphabet size.
 - This is $\mathcal{O}(n / \log_{\sigma} n)$ words in the word RAM model.
 - For example, if $\sigma = 2$, then 64 characters (bits) can be stored in an integer variable.

Is $\mathcal{O}(n)$ space optimal?

- The string T can be encoded in $\mathcal{O}(n \log \sigma)$ bits, where σ is the alphabet size.
 - This is $\mathcal{O}(n / \log_{\sigma} n)$ words in the word RAM model.
 - For example, if $\sigma = 2$, then 64 characters (bits) can be stored in an integer variable.
- The arrays rank and LCP take approximately $n \log n$ bits each.

Is $\mathcal{O}(n)$ space optimal?

- The string T can be encoded in $\mathcal{O}(n \log \sigma)$ bits, where σ is the alphabet size.
 - This is $\mathcal{O}(n / \log_{\sigma} n)$ words in the word RAM model.
 - For example, if $\sigma = 2$, then 64 characters (bits) can be stored in an integer variable.
- The arrays rank and LCP take approximately $n \log n$ bits each.

Questions:

- 1 What query time can be achieved in $\mathcal{O}(n \log \sigma)$ bits of space?

Is $\mathcal{O}(n)$ space optimal?

- The string T can be encoded in $\mathcal{O}(n \log \sigma)$ bits, where σ is the alphabet size.
 - This is $\mathcal{O}(n / \log_{\sigma} n)$ words in the word RAM model.
 - For example, if $\sigma = 2$, then 64 characters (bits) can be stored in an integer variable.
- The arrays rank and LCP take approximately $n \log n$ bits each.

Questions:

- 1 What query time can be achieved in $\mathcal{O}(n \log \sigma)$ bits of space?
 - Still $\mathcal{O}(1)$ query time, with $\mathcal{O}(n / \log_{\sigma} n)$ construction time.

Is $\mathcal{O}(n)$ space optimal?

- The string T can be encoded in $\mathcal{O}(n \log \sigma)$ bits, where σ is the alphabet size.
 - This is $\mathcal{O}(n / \log_{\sigma} n)$ words in the word RAM model.
 - For example, if $\sigma = 2$, then 64 characters (bits) can be stored in an integer variable.
- The arrays rank and LCP take approximately $n \log n$ bits each.

Questions:

- 1 What query time can be achieved in $\mathcal{O}(n \log \sigma)$ bits of space?
 - Still $\mathcal{O}(1)$ query time, with $\mathcal{O}(n / \log_{\sigma} n)$ construction time.
- 2 What about s extra space on top of T for $1 \leq s \leq n / \log_{\sigma} n$?

Is $\mathcal{O}(n)$ space optimal?

- The string T can be encoded in $\mathcal{O}(n \log \sigma)$ bits, where σ is the alphabet size.
 - This is $\mathcal{O}(n / \log_{\sigma} n)$ words in the word RAM model.
 - For example, if $\sigma = 2$, then 64 characters (bits) can be stored in an integer variable.
- The arrays rank and LCP take approximately $n \log n$ bits each.

Questions:

- 1 What query time can be achieved in $\mathcal{O}(n \log \sigma)$ bits of space?
 - Still $\mathcal{O}(1)$ query time, with $\mathcal{O}(n / \log_{\sigma} n)$ construction time.
- 2 What about s extra space on top of T for $1 \leq s \leq n / \log_{\sigma} n$?
 - $\mathcal{O}(n / (s \log_{\sigma} n))$ query time!

- A naive algorithm computes $\ell = \text{LCE}(i, j)$ in $\mathcal{O}(1 + \ell)$ time (no data structure needed).

- A naive algorithm computes $\ell = \text{LCE}(i, j)$ in $\mathcal{O}(1 + \ell)$ time (no data structure needed).
- If the machine word fits α characters, this can be improved to $\mathcal{O}(1 + \ell/\alpha)$ time.
 - This is always $\mathcal{O}(1 + \ell/\log_\sigma n)$ time because $\alpha = \Omega(\log_\sigma n)$.

Exploiting Bit-Parallelism

- A naive algorithm computes $\ell = \text{LCE}(i, j)$ in $\mathcal{O}(1 + \ell)$ time (no data structure needed).
- If the machine word fits α characters, this can be improved to $\mathcal{O}(1 + \ell/\alpha)$ time.
 - This is always $\mathcal{O}(1 + \ell/\log_\sigma n)$ time because $\alpha = \Omega(\log_\sigma n)$.

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

LCE(2, 9) =

Exploiting Bit-Parallelism

- A naive algorithm computes $\ell = \text{LCE}(i, j)$ in $\mathcal{O}(1 + \ell)$ time (no data structure needed).
- If the machine word fits α characters, this can be improved to $\mathcal{O}(1 + \ell/\alpha)$ time.
 - This is always $\mathcal{O}(1 + \ell/\log_{\sigma} n)$ time because $\alpha = \Omega(\log_{\sigma} n)$.

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

$\text{LCE}(2, 9) =$

Extract $T[i..i + \alpha]$: 1 0 1 0 0 1 0 1

Exploiting Bit-Parallelism

- A naive algorithm computes $\ell = \text{LCE}(i, j)$ in $\mathcal{O}(1 + \ell)$ time (no data structure needed).
- If the machine word fits α characters, this can be improved to $\mathcal{O}(1 + \ell/\alpha)$ time.
 - This is always $\mathcal{O}(1 + \ell/\log_{\sigma} n)$ time because $\alpha = \Omega(\log_{\sigma} n)$.

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

$\text{LCE}(2, 9) =$

Extract $T[i..i + \alpha)$: 1 0 1 0 0 1 0 1
Extract $T[j..j + \alpha)$: 1 0 1 0 1 0 0 1

Exploiting Bit-Parallelism

- A naive algorithm computes $\ell = \text{LCE}(i, j)$ in $\mathcal{O}(1 + \ell)$ time (no data structure needed).
- If the machine word fits α characters, this can be improved to $\mathcal{O}(1 + \ell/\alpha)$ time.
 - This is always $\mathcal{O}(1 + \ell/\log_{\sigma} n)$ time because $\alpha = \Omega(\log_{\sigma} n)$.

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

LCE(2, 9) =

```
Extract  $T[i..i + \alpha)$ : 1 0 1 0 0 1 0 1
Extract  $T[j..j + \alpha)$ : 1 0 1 0 1 0 0 1
xor: 0 0 0 0 1 1 0 0
```

Exploiting Bit-Parallelism

- A naive algorithm computes $\ell = \text{LCE}(i, j)$ in $\mathcal{O}(1 + \ell)$ time (no data structure needed).
- If the machine word fits α characters, this can be improved to $\mathcal{O}(1 + \ell/\alpha)$ time.
 - This is always $\mathcal{O}(1 + \ell/\log_\sigma n)$ time because $\alpha = \Omega(\log_\sigma n)$.

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

$$\text{LCE}(2, 9) = 4$$

Extract $T[i..i + \alpha)$:	1 0 1 0 0 1 0 1
Extract $T[j..j + \alpha)$:	1 0 1 0 1 0 0 1
xor:	0 0 0 0 1 1 0 0
clz:	4

LCE Queries with a String Synchronizing Set

- 1 Build the rank and LCP tables only for suffixes starting at *synchronizing positions* (set S).

1 1 0 1 0 0 1 0 1 0 1 0 1 0 0 1 0 1 0 0

LCE Queries with a String Synchronizing Set

- 1 Build the rank and LCP tables only for suffixes starting at *synchronizing positions* (set S).

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

-	1	010100
6	2	01010010100
5	3	0101010010100
7	4	010101010010100
0	5	100
3	6	10010100
7	7	10010101010010100
2	8	101010010100
6	9	10101010010100
1	10	11010010101010010100

LCE Queries with a String Synchronizing Set

- 1 Build the rank and LCP tables only for suffixes starting at *synchronizing positions* (set S).
- 2 Reduce arbitrary $LCE(i, j)$ to $LCE(i + \delta, j + \delta)$ for $i + \delta, j + \delta \in S$.

$\frac{1}{10} \frac{2}{1} \frac{3}{0} \frac{4}{7} \frac{5}{0} \frac{6}{4} \frac{7}{9} \frac{8}{3} \frac{9}{8} \frac{10}{2} \frac{11}{1} \frac{12}{0} \frac{13}{6} \frac{14}{0} \frac{15}{1} \frac{16}{1} \frac{17}{0} \frac{18}{5} \frac{19}{0} \frac{20}{0}$

$LCE(16, 7) =$

-	1	010100
6	2	01010010100
5	3	0101010010100
7	4	010101010010100
0	5	100
3	6	10010100
7	7	10010101010010100
2	8	101010010100
6	9	10101010010100
1	10	11010010101010010100

LCE Queries with a String Synchronizing Set

- 1 Build the rank and LCP tables only for suffixes starting at *synchronizing positions* (set S).
- 2 Reduce arbitrary $LCE(i, j)$ to $LCE(i + \delta, j + \delta)$ for $i + \delta, j + \delta \in S$.

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

$$LCE(16, 7) = LCE(18, 9) + 2 =$$

-	1	010100
6	2	01010010100
5	3	0101010010100
7	4	010101010010100
0	5	100
3	6	10010100
7	7	10010101010010100
2	8	101010010100
6	9	10101010010100
1	10	11010010101010010100

LCE Queries with a String Synchronizing Set

- 1 Build the rank and LCP tables only for suffixes starting at *synchronizing positions* (set S).
- 2 Reduce arbitrary $LCE(i, j)$ to $LCE(i + \delta, j + \delta)$ for $i + \delta, j + \delta \in S$.

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

$$LCE(16, 7) = LCE(18, 9) + 2 =$$

-	1	010100
6	2	01010010100
5	3	0101010010100
7	4	010101010010100
0	5	100
3	6	10010100
7	7	10010101010010100
2	8	101010010100
6	9	10101010010100
1	10	11010010101010010100

LCE Queries with a String Synchronizing Set

- 1 Build the rank and LCP tables only for suffixes starting at *synchronizing positions* (set S).
- 2 Reduce arbitrary $LCE(i, j)$ to $LCE(i + \delta, j + \delta)$ for $i + \delta, j + \delta \in S$.

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

$$LCE(16, 7) = LCE(18, 9) + 2 = 2 + 2 = 4$$

-	1	010100
6	2	01010010100
5	3	0101010010100
7	4	010101010010100
0	5	100
3	6	10010100
7	7	10010101010010100
2	8	101010010100
6	9	10101010010100
1	10	11010010101010010100

LCE Queries with a String Synchronizing Set

- 1 Build the rank and LCP tables only for suffixes starting at *synchronizing positions* (set S).
- 2 Reduce arbitrary $LCE(i, j)$ to $LCE(i + \delta, j + \delta)$ for $i + \delta, j + \delta \in S$.
- 3 Make sure that $LCE(i, j) = \mathcal{O}(n/|S|)$ or there is $\delta = \mathcal{O}(n/|S|)$ with $i + \delta, j + \delta \in S$.
 - This would guarantee runtime $\mathcal{O}(1 + n/(|S| \log_\sigma n))$.

$\frac{1}{10} \frac{1}{7} \frac{0}{7} \frac{1}{4} \frac{0}{4} \frac{0}{9} \frac{1}{3} \frac{0}{8} \frac{1}{2} \frac{0}{2} \frac{1}{6} \frac{0}{6} \frac{0}{1} \frac{1}{5} \frac{0}{5} \frac{1}{5} \frac{0}{5} \frac{0}{5}$

$$LCE(16, 7) = LCE(18, 9) + 2 = 2 + 2 = 4$$

-	1	010100
6	2	01010010100
5	3	0101010010100
7	4	010101010010100
0	5	100
3	6	10010100
7	7	10010101010010100
2	8	101010010100
6	9	10101010010100
1	10	11010010101010010100

LCE Queries with a String Synchronizing Set

- Build the rank and LCP tables only for suffixes starting at *synchronizing positions* (set S).
- Reduce arbitrary $LCE(i, j)$ to $LCE(i + \delta, j + \delta)$ for $i + \delta, j + \delta \in S$.
- Make sure that $LCE(i, j) = \mathcal{O}(n/|S|)$ or there is $\delta = \mathcal{O}(n/|S|)$ with $i + \delta, j + \delta \in S$.
 - This would guarantee runtime $\mathcal{O}(1 + n/(|S| \log_\sigma n))$.

$\frac{1}{10} \frac{2}{1} \frac{3}{0} \frac{4}{7} \frac{5}{0} \frac{6}{4} \frac{7}{9} \frac{8}{3} \frac{9}{8} \frac{10}{2} \frac{11}{1} \frac{12}{0} \frac{13}{6} \frac{14}{0} \frac{15}{1} \frac{16}{5} \frac{17}{0} \frac{18}{1} \frac{19}{0} \frac{20}{0}$

$$LCE(16, 7) = LCE(18, 9) + 2 = 2 + 2 = 4$$

Desired properties of a τ -synchronizing set S :

consistency $T[i..i + 2\tau)$ determines if $i \in S$.

density $S \cap [i..i + \tau) = \emptyset$ for every i .

small size $|S| = \mathcal{O}(n/\tau)$.

-	1	010100
6	2	01010010100
5	3	0101010010100
7	4	010101010010100
0	5	100
3	6	10010100
7	7	10010101010010100
8	8	101010010100
6	9	10101010010100
1	10	11010010101010010100

Desired properties of a τ -synchronizing set S :

consistency $T[i..i + 2\tau)$ determines if $i \in S$.

density $S \cap [i..i + \tau) = \emptyset$ for every i .

small size $|S| = \mathcal{O}(n/\tau)$.

Issues with Periodic Regions

Desired properties of a τ -synchronizing set S :

consistency $T[i..i + 2\tau)$ determines if $i \in S$.

density $S \cap [i..i + \tau) = \emptyset$ for every i .

small size $|S| = \mathcal{O}(n/\tau)$.

Issue: Such a set S cannot always exist.

0 0

Issues with Periodic Regions

Desired properties of a τ -synchronizing set S :

consistency $T[i..i+2\tau)$ determines if $i \in S$.

density $S \cap [i..i+\tau) = \emptyset$ for every i .

small size $|S| = \mathcal{O}(n/\tau)$.

Issue: Such a set S cannot always exist.

0 0

The problematic case is when T contains length- τ substrings with period $o(\tau)$.

Workarounds:

String Synchronizing Sets: Construction

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

String Synchronizing Sets: Construction

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

- 1 Draw a random linear order on the set of length- τ substrings of T :

100 \prec 110 \prec 101 \prec 001 \prec 010

String Synchronizing Sets: Construction

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

- 1 Draw a random linear order on the set of length- τ substrings of T :

100 \prec 110 \prec 101 \prec 001 \prec 010

- 2 Add i to S if the earliest length- τ substring of $T[i..i+2\tau)$ is
 - the prefix $T[i..i+\tau)$, or
 - the suffix $T[i+\tau..i+2\tau)$.

String Synchronizing Sets: Construction

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

- 1 Draw a random linear order on the set of length- τ substrings of T :

$100 \prec 110 \prec 101 \prec 001 \prec 010$

- 2 Add i to S if the earliest length- τ substring of $T[i..i+2\tau)$ is
 - the prefix $T[i..i+\tau)$, or
 - the suffix $T[i+\tau..i+2\tau)$.

String Synchronizing Sets: Construction

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

- 1 Draw a random linear order on the set of length- τ substrings of T :

$100 \prec 110 \prec 101 \prec 001 \prec 010$

- 2 Add i to S if the earliest length- τ substring of $T[i..i+2\tau)$ is
 - the prefix $T[i..i+\tau)$, or
 - the suffix $T[i+\tau..i+2\tau)$.

String Synchronizing Sets: Construction

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

- 1 Draw a random linear order on the set of length- τ substrings of T :

$100 \prec 110 \prec 101 \prec 001 \prec 010$

- 2 Add i to S if the earliest length- τ substring of $T[i..i+2\tau)$ is
 - the prefix $T[i..i+\tau)$, or
 - the suffix $T[i+\tau..i+2\tau)$.

String Synchronizing Sets: Construction

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

- 1 Draw a random linear order on the set of length- τ substrings of T :

100 \prec 110 \prec 101 \prec 001 \prec 010

- 2 Add i to S if the earliest length- τ substring of $T[i..i+2\tau)$ is
 - the prefix $T[i..i+\tau)$, or
 - the suffix $T[i+\tau..i+2\tau)$.

String Synchronizing Sets: Construction

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

- 1 Draw a random linear order on the set of length- τ substrings of T :

$100 \prec 110 \prec 101 \prec 001 \prec 010$

- 2 Add i to S if the earliest length- τ substring of $T[i..i+2\tau)$ is
 - the prefix $T[i..i+\tau)$, or
 - the suffix $T[i+\tau..i+2\tau)$.

String Synchronizing Sets: Construction

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

- 1 Draw a random linear order on the set of length- τ substrings of T :

$100 \prec 110 \prec 101 \prec 001 \prec 010$

- 2 Add i to S if the earliest length- τ substring of $T[i..i+2\tau)$ is
 - the prefix $T[i..i+\tau)$, or
 - the suffix $T[i+\tau..i+2\tau)$.

String Synchronizing Sets: Construction

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

- 1 Draw a random linear order on the set of length- τ substrings of T :

100 \prec 110 \prec 101 \prec 001 \prec 010

- 2 Add i to S if the earliest length- τ substring of $T[i..i+2\tau)$ is
 - the prefix $T[i..i+\tau)$, or
 - the suffix $T[i+\tau..i+2\tau)$.

String Synchronizing Sets: Construction

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

- 1 Draw a random linear order on the set of length- τ substrings of T :

100 \prec 110 \prec 101 \prec 001 \prec 010

- 2 Add i to S if the earliest length- τ substring of $T[i..i+2\tau)$ is
 - the prefix $T[i..i+\tau)$, or
 - the suffix $T[i+\tau..i+2\tau)$.

String Synchronizing Sets: Construction

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

- 1 Draw a random linear order on the set of length- τ substrings of T :

100 \prec 110 \prec 101 \prec 001 \prec 010

- 2 Add i to S if the earliest length- τ substring of $T[i..i+2\tau)$ is
 - the prefix $T[i..i+\tau)$, or
 - the suffix $T[i+\tau..i+2\tau)$.

String Synchronizing Sets: Construction

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

- 1 Draw a random linear order on the set of length- τ substrings of T :

$100 \prec 110 \prec 101 \prec 001 \prec 010$

- 2 Add i to S if the earliest length- τ substring of $T[i..i+2\tau)$ is
 - the prefix $T[i..i+\tau)$, or
 - the suffix $T[i+\tau..i+2\tau)$.
- 3 Consistency and density are easy to argue.

String Synchronizing Sets: Construction

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

- 1 Draw a random linear order on the set of length- τ substrings of T :

$100 \prec 110 \prec 101 \prec 001 \prec 010$

- 2 Add i to S if the earliest length- τ substring of $T[i..i+2\tau)$ is
 - the prefix $T[i..i+\tau)$, or
 - the suffix $T[i+\tau..i+2\tau)$.
- 3 Consistency and density are easy to argue.
- 4 If no length- τ substrings of $T[i..i+2\tau)$ has period $o(\tau)$, then $\Pr[i \in S] = \mathcal{O}(1/\tau)$.
 - $T[i..i+2\tau)$ contains $\Omega(\tau)$ distinct length- τ substrings.

LCE queries in small space:

- A τ -synchronizing S set with $\tau = \Theta(n/s)$.
- Expected size $\text{Exp}[|S|] = \mathcal{O}(n/\tau)$.
 - Requires adaptations for periodic regions or assuming no length- τ substring has period $o(\tau)$.

LCE queries in small space:

- A τ -synchronizing S set with $\tau = \Theta(n/s)$.
- Expected size $\text{Exp}[|S|] = \mathcal{O}(n/\tau)$.
 - Requires adaptations for periodic regions or assuming no length- τ substring has period $o(\tau)$.
- Worst-case size $\mathcal{O}(n/\tau)$ after $\mathcal{O}(1)$ attempts (in expectation).

LCE queries in small space:

- A τ -synchronizing S set with $\tau = \Theta(n/s)$.
- Expected size $\text{Exp}[|S|] = \mathcal{O}(n/\tau)$.
 - Requires adaptations for periodic regions or assuming no length- τ substring has period $o(\tau)$.
- Worst-case size $\mathcal{O}(n/\tau)$ after $\mathcal{O}(1)$ attempts (in expectation).
- The construction can be derandomized using the method of pessimistic estimators.

LCE queries in small space:

- A τ -synchronizing S set with $\tau = \Theta(n/s)$.
- Expected size $\text{Exp}[|S|] = \mathcal{O}(n/\tau)$.
 - Requires adaptations for periodic regions or assuming no length- τ substring has period $o(\tau)$.
- Worst-case size $\mathcal{O}(n/\tau)$ after $\mathcal{O}(1)$ attempts (in expectation).
- The construction can be derandomized using the method of pessimistic estimators.

Final data structure:

size $\mathcal{O}(|S|) = \mathcal{O}(n/\tau)$
query time $\mathcal{O}(1 + \tau/\log_\sigma n)$

for $\tau = \Theta(\log_\sigma n)$
 $\mathcal{O}(n/\log_\sigma n)$
 $\mathcal{O}(1)$

LCE queries in small space:

- A τ -synchronizing S set with $\tau = \Theta(n/s)$.
- Expected size $\text{Exp}[|S|] = \mathcal{O}(n/\tau)$.
 - Requires adaptations for periodic regions or assuming no length- τ substring has period $o(\tau)$.
- Worst-case size $\mathcal{O}(n/\tau)$ after $\mathcal{O}(1)$ attempts (in expectation).
- The construction can be derandomized using the method of pessimistic estimators.

Final data structure:

size $\mathcal{O}(|S|) = \mathcal{O}(n/\tau)$
query time $\mathcal{O}(1 + \tau/\log_\sigma n)$
construction time $\mathcal{O}(n)$

for $\tau = \Theta(\log_\sigma n)$
 $\mathcal{O}(n/\log_\sigma n)$
 $\mathcal{O}(1)$
 $\mathcal{O}(n/\log_\sigma n)$

LCE queries in small space:

- A τ -synchronizing S set with $\tau = \Theta(n/s)$.
- Expected size $\text{Exp}[|S|] = \mathcal{O}(n/\tau)$.
 - Requires adaptations for periodic regions or assuming no length- τ substring has period $o(\tau)$.
- Worst-case size $\mathcal{O}(n/\tau)$ after $\mathcal{O}(1)$ attempts (in expectation).
- The construction can be derandomized using the method of pessimistic estimators.

Final data structure:

size $\mathcal{O}(|S|) = \mathcal{O}(n/\tau)$
query time $\mathcal{O}(1 + \tau/\log_\sigma n)$
construction time $\mathcal{O}(n)$

for $\tau = \Theta(\log_\sigma n)$
 $\mathcal{O}(n/\log_\sigma n)$
 $\mathcal{O}(1)$
 $\mathcal{O}(n/\log_\sigma n)$

Open question:

- $\mathcal{O}(n/\log_\sigma n)$ construction time for $\tau = \Omega(\log_\sigma n)$.

Further applications of synchronizing sets:

KRRW, SODA'15 Internal Pattern Matching queries, Period queries

KK, STOC'19 Burrows–Wheeler Transform construction (bzip2, etc.)

CKPR, ESA'21 Longest Common Substring problem

Further applications of synchronizing sets:

- KRRW, SODA'15** Internal Pattern Matching queries, Period queries
- KK, STOC'19** Burrows–Wheeler Transform construction (bzip2, etc.)
- CKPR, ESA'21** Longest Common Substring problem
- KK, STOC'22** dynamic suffix array

Further applications of synchronizing sets:

- KRRW, SODA'15** Internal Pattern Matching queries, Period queries
- KK, STOC'19** Burrows–Wheeler Transform construction (bzip2, etc.)
- CKPR, ESA'21** Longest Common Substring problem
- KK, STOC'22** dynamic suffix array
- AJ, SODA'22** quantum algorithm for Longest Common Substring
- JN, SODA'23** quantum data structure for LCE queries

Further applications of synchronizing sets:

- KRRW, SODA'15** Internal Pattern Matching queries, Period queries
- KK, STOC'19** Burrows–Wheeler Transform construction (bzip2, etc.)
- CKPR, ESA'21** Longest Common Substring problem
- KK, STOC'22** dynamic suffix array
- AJ, SODA'22** quantum algorithm for Longest Common Substring
- JN, SODA'23** quantum data structure for LCE queries
- KK, SODA'23** compressed suffix arrays and text indexes

Further applications of synchronizing sets:

- KRRW, SODA'15** Internal Pattern Matching queries, Period queries
- KK, STOC'19** Burrows–Wheeler Transform construction (bzip2, etc.)
- CKPR, ESA'21** Longest Common Substring problem
- KK, STOC'22** dynamic suffix array
- AJ, SODA'22** quantum algorithm for Longest Common Substring
- JN, SODA'23** quantum data structure for LCE queries
- KK, SODA'23** compressed suffix arrays and text indexes

Implementations:

- DFHKK, ESA'20** LCE queries, ignoring the periodic case
- in progress** LCE queries, with worst-case guarantees

Further applications of synchronizing sets:

- KRRW, SODA'15** Internal Pattern Matching queries, Period queries
- KK, STOC'19** Burrows–Wheeler Transform construction (bzip2, etc.)
- CKPR, ESA'21** Longest Common Substring problem
- KK, STOC'22** dynamic suffix array
- AJ, SODA'22** quantum algorithm for Longest Common Substring
- JN, SODA'23** quantum data structure for LCE queries
- KK, SODA'23** compressed suffix arrays and text indexes

Implementations:

- DFHKK, ESA'20** LCE queries, ignoring the periodic case
- in progress** LCE queries, with worst-case guarantees

Related techniques in applications (w/o theoretical guarantees):

- Bioinformatics (minimizers)
- Plagiarism detection (winnowing)
- Storage systems (content-based chunking)