

Parallel Computation of Matching Statistics and Average Common Substring

Fabio Garofalo Daniele Greco²
Giovanna Rosone² **Marinella Sciortino**¹

¹Dipartimento di Matematica e Informatica
Università di Palermo, Italy

²Dipartimento di Informatica
Università di Pisa, Italy

BITS 2019 - Analysis of Big Omics Data
June 26th, 2019

Motivation

- Increased availability of large sets of biological sequences
- Tools for sequence comparison *a fortiori* need alignment-free based approaches
- Most alignment-free approaches require the computation of statistics when comparing sequences
- Such computations may not scale well in internal memory when very large collections of long sequences are considered

Motivation

- Increased availability of large sets of biological sequences
- Tools for sequence comparison *a fortiori* need alignment-free based approaches
- Most alignment-free approaches require the computation of statistics when comparing sequences
- Such computations may not scale well in internal memory when very large collections of long sequences are considered
- **Our contribution:**

Motivation

- Increased availability of large sets of biological sequences
- Tools for sequence comparison *a fortiori* need alignment-free based approaches
- Most alignment-free approaches require the computation of statistics when comparing sequences
- Such computations may not scale well in internal memory when very large collections of long sequences are considered
- **Our contribution:**
 - the **Colored Longest Common Prefix** array: a new conceptual data structure that uses an alignment-free approach and can be computed via sequential scans in semi-external memory

Motivation

- Increased availability of large sets of biological sequences
- Tools for sequence comparison *a fortiori* need alignment-free based approaches
- Most alignment-free approaches require the computation of statistics when comparing sequences
- Such computations may not scale well in internal memory when very large collections of long sequences are considered
- **Our contribution:**
 - the **Colored Longest Common Prefix** array: a new conceptual data structure that uses an alignment-free approach and can be computed via sequential scans in semi-external memory
 - it implicitly stores all the information necessary **to compute statistics** on distinguishing, repeating, or matching substrings **within collections of strings**.

Motivation

- Increased availability of large sets of biological sequences
- Tools for sequence comparison *a fortiori* need alignment-free based approaches
- Most alignment-free approaches require the computation of statistics when comparing sequences
- Such computations may not scale well in internal memory when very large collections of long sequences are considered
- **Our contribution:**
 - the **Colored Longest Common Prefix** array: a new conceptual data structure that uses an alignment-free approach and can be computed via sequential scans in semi-external memory
 - it implicitly stores all the information necessary **to compute statistics** on distinguishing, repeating, or matching substrings **within collections of strings**.
 - efficient lightweight strategy to solve the **multi-string Average Common Substring (ACS) Problem** and experimental results.

Preliminaries

- $\Sigma = \{c_1, c_2, \dots, c_\sigma\}$ be a finite ordered alphabet
- \mathcal{S} is a collection of m strings over Σ
- n_i is the length of the string s_i
- A distinct end-marker symbol $\$i < c_1$ is appended to each string s_i
- $N = \sum_{i=1}^m n_i + m$ is the length of the collection \mathcal{S}
- Each string (or subset of strings) is identified by a specific color

EBWT (Mantaci, Restivo, Rosone, S. 2007)

Let $\mathcal{S} = \{s_1, s_2, \dots, s_m\}$ be a collection of strings.

String collection \mathcal{S}

	0	1	2	3	4	5	6
s_1	G	C	C	A	A	C	$\$1$
s_2	G	A	G	C	T	C	$\$2$
s_3	T	C	G	C	T	T	$\$3$

EBWT (Mantaci, Restivo, Rosone, S. 2007)

Let $\mathcal{S} = \{s_1, s_2, \dots, s_m\}$ be a collection of strings.

- The extended Burrows-Wheeler Transform for a string collection \mathcal{S} , known as EBWT or multi-string BWT, is a reversible transformation that produces a string (denoted by $\text{ebwt}(\mathcal{S})$) that is a permutation of the characters of all strings in \mathcal{S}

		String collection \mathcal{S}					
	0	1	2	3	4	5	6
s_1	G	C	C	A	A	C	$\$1$
s_2	G	A	G	C	T	C	$\$2$
s_3	T	C	G	C	T	T	$\$3$

EBWT (Mantaci, Restivo, Rosone, S. 2007)

Let $\mathcal{S} = \{s_1, s_2, \dots, s_m\}$ be a collection of strings.

- The extended Burrows-Wheeler Transform for a string collection \mathcal{S} , known as EBWT or multi-string BWT, is a reversible transformation that produces a string (denoted by $\text{ebwt}(\mathcal{S})$) that is a permutation of the characters of all strings in \mathcal{S}

		String collection \mathcal{S}						
		0	1	2	3	4	5	6
s_1		G	C	C	A	A	C	$\$1$
s_2		G	A	G	C	T	C	$\$2$
s_3		T	C	G	C	T	T	$\$3$

$\text{ebwt}(\mathcal{S})$	Sorted Suffixes of \mathcal{S}
C	$\$1$
C	$\$2$
T	$\$3$
C	AAC $\$1$
A	AC $\$1$
G	AGCTC $\$2$
A	C $\$1$
T	C $\$2$
C	CAAC $\$1$
G	CCAAC $\$1$
T	CGCTT $\$3$
G	CTC $\$2$
G	CTT $\$3$
$\$2$	GAGCTC $\$2$
$\$1$	GCCAAC $\$1$
A	GCTC $\$2$
C	GCTT $\$3$
T	T $\$3$
C	TC $\$2$
$\$3$	TCGCTT $\$3$
C	TT $\$3$

- Sort all the suffixes of the strings in \mathcal{S} ;

EBWT (Mantaci, Restivo, Rosone, S. 2007)

Let $\mathcal{S} = \{s_1, s_2, \dots, s_m\}$ be a collection of strings.

- The extended Burrows-Wheeler Transform for a string collection \mathcal{S} , known as EBWT or multi-string BWT, is a reversible transformation that produces a string (denoted by $\text{ebwt}(\mathcal{S})$) that is a permutation of the characters of all strings in \mathcal{S}

		String collection \mathcal{S}						
		0	1	2	3	4	5	6
s_1		G	C	C	A	A	C	$\$1$
s_2		G	A	G	C	T	C	$\$2$
s_3		T	C	G	C	T	T	$\$3$

$\text{ebwt}(\mathcal{S})$	Sorted Suffixes of \mathcal{S}
C	$\$1$
C	$\$2$
T	$\$3$
C	AAC $\$1$
A	AC $\$1$
G	AGCTC $\$2$
A	C $\$1$
T	C $\$2$
C	CAAC $\$1$
G	CCAAC $\$1$
T	CGCTT $\$3$
G	CTC $\$2$
G	CTT $\$3$
$\$2$	GAGCTC $\$2$
$\$1$	GCCAAC $\$1$
A	GCTC $\$2$
C	GCTT $\$3$
T	T $\$3$
C	TC $\$2$
$\$3$	TCGCTT $\$3$
C	TT $\$3$

- Sort all the suffixes of the strings in \mathcal{S} ;
- The output $\text{ebwt}(\mathcal{S})$ is obtained by concatenating the symbols that (circularly) precede the first symbol of each suffix in the list of (lexicographically) sorted suffixes of \mathcal{S} .

LCP and Colored EBWT

- The *longest common prefix* (LCP) array of the collection \mathcal{S} is the array $\text{lcp}(\mathcal{S})$ of length $N + 1$, such that $\text{lcp}(\mathcal{S})[i]$, with $2 \leq i \leq N$, is the length of the longest common prefix between the suffixes associated to the positions i and $i - 1$ in $\text{ebwt}(\mathcal{S})$. By default, $\text{lcp}(\mathcal{S})[1] = \text{lcp}(\mathcal{S})[N + 1] = -1$

$\text{lcp}(\mathcal{S})$	$\text{id}(\mathcal{S})$	$\text{ebwt}(\mathcal{S})$	Sorted Suffixes of \mathcal{S}
-1	1	C	$\$1$
0	2	C	$\$2$
0	3	T	$\$3$
0	1	C	AAC $\$1$
1	1	A	AC $\$1$
1	2	G	AGCTC $\$2$
0	1	A	C $\$1$
1	2	T	C $\$2$
1	1	C	CAAC $\$1$
1	1	G	CCAAC $\$1$
1	3	T	CGCTT $\$3$
1	2	G	CTC $\$2$
2	3	G	CTT $\$3$
0	2	$\$2$	GAGCTC $\$2$
1	1	$\$1$	GCCAAC $\$1$
2	2	A	GCTC $\$2$
3	3	C	GCTT $\$3$
0	3	T	T $\$3$
1	2	C	TC $\$2$
2	3	$\$3$	TCGCTT $\$3$
1	3	C	TT $\$3$
-1			

LCP and Colored EBWT

- The *longest common prefix* (LCP) array of the collection \mathcal{S} is the array $\text{lcp}(\mathcal{S})$ of length $N + 1$, such that $\text{lcp}(\mathcal{S})[i]$, with $2 \leq i \leq N$, is the length of the longest common prefix between the suffixes associated to the positions i and $i - 1$ in $\text{ebwt}(\mathcal{S})$. By default, $\text{lcp}(\mathcal{S})[1] = \text{lcp}(\mathcal{S})[N + 1] = -1$
- $\text{lcp}(i, j)$ the length of the LCP between the suffixes at positions i and j , i.e. $\min\{\text{lcp}(\mathcal{S})[l] : i < l \leq j\}$.

$\text{lcp}(\mathcal{S})$	$\text{id}(\mathcal{S})$	$\text{ebwt}(\mathcal{S})$	Sorted Suffixes of \mathcal{S}
-1	1	C	$\$1$
0	2	C	$\$2$
0	3	T	$\$3$
0	1	C	AAC $\$1$
1	1	A	AC $\$1$
1	2	G	AGCTC $\$2$
0	1	A	C $\$1$
1	2	T	C $\$2$
1	1	C	CAAC $\$1$
1	1	G	CCAAC $\$1$
1	3	T	CGCTT $\$3$
1	2	G	CTC $\$2$
2	3	G	CTT $\$3$
0	2	$\$2$	GAGCTC $\$2$
1	1	$\$1$	GCCAAC $\$1$
2	2	A	GCTC $\$2$
3	3	C	GCTT $\$3$
0	3	T	T $\$3$
1	2	C	TC $\$2$
2	3	$\$3$	TCGCTT $\$3$
1	3	C	TT $\$3$
-1			

LCP and Colored EBWT

- The *longest common prefix* (LCP) array of the collection \mathcal{S} is the array $\text{lcp}(\mathcal{S})$ of length $N + 1$, such that $\text{lcp}(\mathcal{S})[i]$, with $2 \leq i \leq N$, is the length of the longest common prefix between the suffixes associated to the positions i and $i - 1$ in $\text{ebwt}(\mathcal{S})$. By default, $\text{lcp}(\mathcal{S})[1] = \text{lcp}(\mathcal{S})[N + 1] = -1$
- $\text{lcp}(i, j)$ the length of the LCP between the suffixes at positions i and j , i.e. $\min\{\text{lcp}(\mathcal{S})[l] : i < l \leq j\}$.
- The output string $\text{ebwt}(\mathcal{S})$, enhanced with the N-integer array of colors $\text{id}(\mathcal{S})$ where $\text{id}(\mathcal{S})[i] = r$, with $1 \leq r \leq m$ and $1 \leq i \leq N$, if $\text{ebwt}(\mathcal{S})[i]$ is a symbol of the string $s_r \in \mathcal{S}$, is called *colored* EBWT.

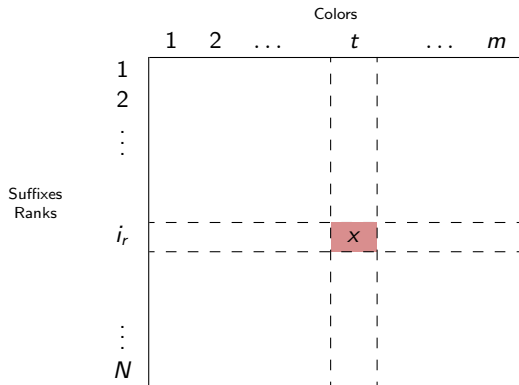
$\text{lcp}(\mathcal{S})$	$\text{id}(\mathcal{S})$	$\text{ebwt}(\mathcal{S})$	Sorted Suffixes of \mathcal{S}
-1	1	C	$\$1$
0	2	C	$\$2$
0	3	T	$\$3$
0	1	C	AAC $\$1$
1	1	A	AC $\$1$
1	2	G	AGCTC $\$2$
0	1	A	C $\$1$
1	2	T	C $\$2$
1	1	C	CAAC $\$1$
1	1	G	CCAAC $\$1$
1	3	T	CGCTT $\$3$
1	2	G	CTC $\$2$
2	3	G	CTT $\$3$
0	2	$\$2$	GAGCTC $\$2$
1	1	$\$1$	GCCAAC $\$1$
2	2	A	GCTC $\$2$
3	3	C	GCTT $\$3$
0	3	T	T $\$3$
1	2	C	TC $\$2$
2	3	$\$3$	TCGCTT $\$3$
1	3	C	TT $\$3$
-1			

Some recent lightweight implementations of EBWT and LCP

- Bauer, M., Cox, A., Rosone, G.: Lightweight algorithms for constructing and inverting the BWT of string collections. *Theor. Comput. Sci.* 483(0), 134–148 (2013) (known as BCR implementation)
- Cox, A.J., Garofalo, F., Rosone, G., Sciortino, M.: Lightweight LCP construction for very large collections of strings. *J. Discrete Algorithms* 37, 17–33 (2016)
- Egidi, L., Louza, F.A., Manzini, G., Telles, G.P.: External memory BWT and LCP computation for sequence collections with applications. *WABI 2018*
- Louza, F., Telles, G., Hoffmann, S., Ciferri, C.: Generalized enhanced suffix array construction in external memory. *Algorithms Mol. Biol.* 12(1), 26 (2017)

Colored Longest Common Prefix (cLCP) array

cLCP is an $(N \times m)$ -integer array representing the longest common prefix between any specific suffix of a (r -colored) string $s_r \in \mathcal{S}$ and the nearest suffixes of a specific (t -colored) string $s_t \in \mathcal{S}$ in the sorted list of suffixes of \mathcal{S} .



$x = \text{LCP}$ value between the r -colored suffix of rank i_r and the nearest t -colored suffix in the sorted list of suffixes

Colored LCP - Formal Definition

Given $1 \leq i_r \leq N$ and $t = 1, \dots, m$,
how $\text{cLCP}[i_r][t]$ is defined?

$$\text{prev}(i, t) = \max\{x \mid 1 \leq x < i, \text{id}(S)[x] = t\}$$

$$\text{next}(i, t) = \min\{x \mid i < x \leq N, \text{id}(S)[x] = t\}$$

i	$\text{lcp}(S)$	$\text{id}(S)$
\vdots	\vdots	\vdots
$\text{prev}(i_r, t)$	l_1	t
\vdots	\vdots	\uparrow $\neq t$ \downarrow
i_r	l	r
\vdots	\vdots	\uparrow $\neq t$ \downarrow
$\text{next}(i_r, t)$	l_2	t
\vdots	\vdots	\vdots

Colored LCP - Formal Definition

Given $1 \leq i_r \leq N$ and $t = 1, \dots, m$,
 how $\text{cLCP}[i_r][t]$ is defined?

$$\text{prev}(i, t) = \max\{x \mid 1 \leq x < i, \text{id}(\mathcal{S})[x] = t\}$$

$$\text{next}(i, t) = \min\{x \mid i < x \leq N, \text{id}(\mathcal{S})[x] = t\}$$

i	$\text{lcp}(\mathcal{S})$	$\text{id}(\mathcal{S})$
\vdots	\vdots	\vdots
$\text{prev}(i_r, t)$	l_1	t
\vdots	\vdots	\uparrow $\neq t$ \downarrow
i_r	l	r
\vdots	\vdots	\uparrow $\neq t$ \downarrow
$\text{next}(i_r, t)$	l_2	t
\vdots	\vdots	\vdots

Upper Colored LCP

$$\text{UcLCP}[i_r][t] = \text{LCP}(\text{prev}(i_r, t), i_r)$$

Colored LCP - Formal Definition

Given $1 \leq i_r \leq N$ and $t = 1, \dots, m$,
 how $\text{cLCP}[i_r][t]$ is defined?

$$\text{prev}(i, t) = \max\{x \mid 1 \leq x < i, \text{id}(\mathcal{S})[x] = t\}$$

$$\text{next}(i, t) = \min\{x \mid i < x \leq N, \text{id}(\mathcal{S})[x] = t\}$$

i	$\text{lcp}(\mathcal{S})$	$\text{id}(\mathcal{S})$
\vdots	\vdots	\vdots
$\text{prev}(i_r, t)$	l_1	t
\vdots	\vdots	\uparrow $\neq t$ \downarrow
i_r	l	r
\vdots	\vdots	\uparrow $\neq t$ \downarrow
$\text{next}(i_r, t)$	l_2	t
\vdots	\vdots	\vdots

Lower Colored LCP

$$\text{LcLCP}[i_r][t] = \text{LCP}(i_r, \text{next}(i_r, t))$$

Colored LCP - Formal Definition

Given $1 \leq i_r \leq N$ and $t = 1, \dots, m$,
 how $\text{cLCP}[i_r][t]$ is defined?

$$\text{prev}(i, t) = \max\{x \mid 1 \leq x < i, \text{id}(\mathcal{S})[x] = t\}$$

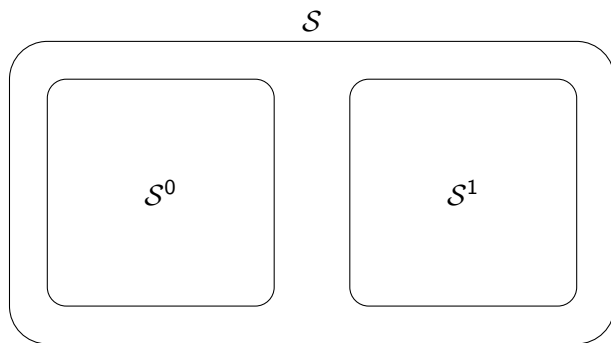
$$\text{next}(i, t) = \min\{x \mid i < x \leq N, \text{id}(\mathcal{S})[x] = t\}$$

i	$\text{lcp}(\mathcal{S})$	$\text{id}(\mathcal{S})$
\vdots	\vdots	\vdots
$\text{prev}(i_r, t)$	l_1	t
\vdots	\vdots	\uparrow $\neq t$ \downarrow
i_r	l	r
\vdots	\vdots	\uparrow $\neq t$ \downarrow
$\text{next}(i_r, t)$	l_2	t
\vdots	\vdots	\vdots

$$\text{cLCP}[i_r][t] = \max\{\text{UcLCP}[i_r][t], \text{LcLCP}[i_r][t]\}$$

Colored LCP on disjoint collections

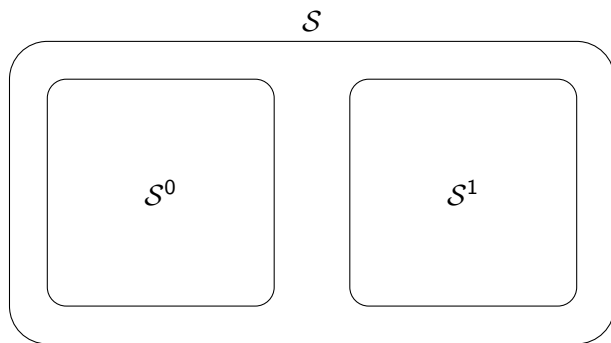
cLCP can also be defined for disjoint collections of strings.



The value $\text{cLCP}[i_r][t]$ is defined for each pair (i_r, t) such that $\text{id}(S)[i_r] = r$, $t \in \text{ID}$ and s_r, s_t belongs to different collections.

Colored LCP on disjoint collections

cLCP can also be defined for disjoint collections of strings.



The value $\text{cLCP}[i_r][t]$ is defined for each pair (i_r, t) such that $\text{id}(\mathcal{S})[i_r] = r$, $t \in \text{ID}$ and s_r, s_t belongs to different collections.

Assume $\mathcal{S}^0 = \{s_\chi\}$ and $\mathcal{S}^1 = \mathcal{S} \setminus \{s_\chi\}$.

χ – intervals

A given string $s_\chi \in \mathcal{S}^0$ implicitly induces a partition of $\text{lcp}(\mathcal{S})$ into open intervals delimited by consecutive suffixes having color χ (or the positions 1 and $N + 1$ of lcp), called χ -intervals. Let us consider a position i_r contained within a χ -interval such that $\text{id}[i_r] = r$ and $s_r \in \mathcal{S}^1$.

i	$\text{lcp}(\mathcal{S})$	$\text{id}(\mathcal{S})$
\vdots	\vdots	\vdots
$\chi_1 = \text{prev}(i_r, \chi)$		χ
		\top
\vdots	\vdots	$\neq \chi$
		\perp
i_r		r
		\top
\vdots	\vdots	$\neq \chi$
		\perp
$\chi_2 = \text{next}(i_r, \chi)$		χ
\vdots	\vdots	\vdots

χ – intervals

A given string $s_\chi \in \mathcal{S}^0$ implicitly induces a partition of $\text{lcp}(\mathcal{S})$ into open intervals delimited by consecutive suffixes having color χ (or the positions 1 and $N + 1$ of lcp), called χ -intervals. Let us consider a position i_r contained within a χ -interval such that $\text{id}[i_r] = r$ and $s_r \in \mathcal{S}^1$.

i	$\text{lcp}(\mathcal{S})$	$\text{id}(\mathcal{S})$
\vdots	\vdots	\vdots
$\chi_1 = \text{prev}(i_r, \chi)$		χ
		\top
\vdots	\vdots	$\neq \chi$
		\perp
i_r		r
		\top
\vdots	\vdots	$\neq \chi$
		\perp
$\chi_2 = \text{next}(i_r, \chi)$		χ
\vdots	\vdots	\vdots

Inside a χ -interval, how to compute $\text{UcLCP}[i_r][\chi]$ and $\text{LcLCP}[i_r][\chi]$ by a lightweight strategy?

UcLCP[i_r][χ] Computation

- ▷ Keep track of the minimum lcp value since the beginning of each χ -interval.

$$\text{UcLCP}[i_r][\chi] = \text{LCP}(\chi_1, i_r) = \min\{\text{lcp}[x] : x \in (\chi_1, i_r]\} = \alpha$$

$\alpha \leftarrow \infty$

i	$\text{lcp}(\mathcal{S})$	$\text{id}(\mathcal{S})$
\vdots	\vdots	\vdots
$\chi_1 = \text{prev}(i_r, \chi)$		χ
\vdots	\vdots	\top $\neq \chi$
i_r		\perp r
\vdots	\vdots	\top $\neq \chi$
$\chi_2 = \text{next}(i_r, \chi)$		\perp χ
\vdots	\vdots	\vdots

UcLCP[i_r][χ] Computation

- ▷ Keep track of the minimum lcp value since the beginning of each χ -interval.

$$\text{UcLCP}[i_r][\chi] = \text{LCP}(\chi_1, i_r) = \min\{\text{lcp}[x] : x \in (\chi_1, i_r]\} = \alpha$$

$\alpha \leftarrow \text{LCP}(\chi_1, i_r)$

i	$\text{lcp}(\mathcal{S})$	$\text{id}(\mathcal{S})$
\vdots	\vdots	\vdots
$\chi_1 = \text{prev}(i_r, \chi)$		χ
\vdots		\top
		$\neq \chi$
		\perp
i_r		r
\vdots		\top
		$\neq \chi$
		\perp
$\chi_2 = \text{next}(i_r, \chi)$		χ
\vdots		\vdots

Example

$$s_X = \text{ACGCGCC}\$X,$$

$$s_1 = \text{ACGAGACGAT}\$1,$$

$$s_2 = \text{AACGCCGCCGGCA}\$2$$

#	id	S	ebwt	lcp	lcp _X	α	UcLCP			LcLCP	cLCP	Sorted suffixes of S
							χ	1	2			
1	χ	1	C	-1	-1							s_X
2	1	0	T	0								s_1
3	2	0	A	0								s_2
4	2	0	C	0								A s_2
5	2	0	s_2	1								A A C G C G C C G G C A s_2
6	1	0	s_1	1								A C G A T s_1
7	1	0	G	4								A C G A T s_1
8	2	0	A	3								A C G C C G C G G C A s_2
9	χ	1	s_0	4	0							A C G C C G C s_X
10	1	0	G	1								A G A C G A T s_1
11	1	0	G	1								A T s_1
12	χ	1	C	0	0							C s_X
13	2	0	G	1								C A s_2
14	χ	1	G	1	1	∞						C C s_X
15	2	0	G	2								C C G C C G C A s_2
16	2	0	G	3								C C G C A s_2
17	1	0	A	1								C G A G A C G A T s_1
18	1	0	A	3								C G A T s_1
19	χ	1	G	2	1							C G C C s_X
20	2	0	A	4								C G C C G C G G C A s_2
21	2	0	C	5								C G C C G G C A s_2
22	χ	1	A	3	3							C G C G C s_X
23	2	0	C	2								C G C A s_2
24	1	0	A	0								G A C G A T s_1
25	1	0	C	2								G A G A C G A T s_1
26	1	0	C	2								G A T s_1
27	2	0	G	1								G C A s_2
28	χ	1	C	2	0							G C C s_X
29	2	0	C	3								G C C G C G G C A s_2
30	2	0	C	4								G C C C G C A s_2
31	χ	1	C	2	2							G C G C s_X
32	2	0	C	1								G C C A s_2
33	1	0	A	0								T s_1
				-1	-1							

Example

$$s_X = \text{ACGCGCC}\$X,$$

$$s_1 = \text{ACGAGACGAT}\$1,$$

$$s_2 = \text{AACGCCGCCGGCA}\$2$$

#	id	S	ebwt	lcp	lcp _X	α	UcLCP			LcLCP	cLCP	Sorted suffixes of S
							χ	1	2			
1	χ	1	C	-1	-1							s_X
2	1	0	T	0								s_1
3	2	0	A	0								s_2
4	2	0	C	0								A s_2
5	2	0	s_2	1								A A C G C G C C G G C A s_2
6	1	0	s_1	1								A C G A T s_1
7	1	0	G	4								A C G A T s_1
8	2	0	A	3								A C G C C G C G G C A s_2
9	χ	1	s_0	4	0							A C G C C G C s_X
10	1	0	G	1								A G A C G A T s_1
11	1	0	G	1								A T s_1
12	χ	1	C	0	0							C s_X
13	2	0	G	1								C A s_2
14	χ	1	G	1	1	∞						C C s_X
15	2	0	G	2	2	2						C C G C C G C A s_2
16	2	0	G	3								C C G G C A s_2
17	1	0	A	1								C G A G A C G A T s_1
18	1	0	A	3								C G A T s_1
19	χ	1	G	2	1							C G C C s_X
20	2	0	A	4								C G C C G C G G C A s_2
21	2	0	C	5								C G C C G G C A s_2
22	χ	1	A	3	3							C G C G C s_X
23	2	0	C	2								C G G C A s_2
24	1	0	A	0								G A C G A T s_1
25	1	0	C	2								G A G A C G A T s_1
26	1	0	C	2								G A T s_1
27	2	0	G	1								G C A s_2
28	χ	1	C	2	0							G C C C s_X
29	2	0	C	3								G C C C G C G G C A s_2
30	2	0	C	4								G C C C G C A s_2
31	χ	1	C	2	2							G C G C s_X
32	2	0	C	1								G C C A s_2
33	1	0	A	0								T s_1
				-1	-1							

Example

$$s_X = \text{ACGCGCC}\$X,$$

$$s_1 = \text{ACGAGACGAT}\$1,$$

$$s_2 = \text{AACGCCGCCGGCA}\$2$$

#	id	S	ebwt	lcp	lcp _X	α	UcLCP			LcLCP	cLCP	Sorted suffixes of S
							χ	1	2			
1	χ	1	C	-1	-1							s_X
2	1	0	T	0								s_1
3	2	0	A	0								s_2
4	2	0	C	0								A s_2
5	2	0	s_2	1								A A C G C G C C G G C A s_2
6	1	0	s_1	1								A C G A T s_1
7	1	0	G	4								A C G A T s_1
8	2	0	A	3								A C G C C G C G G C A s_2
9	χ	1	s_0	4	0							A C G C C G C s_X
10	1	0	G	1								A G A C G A T s_1
11	1	0	G	1								A T s_1
12	χ	1	C	0	0							C s_X
13	2	0	G	1								C A s_2
14	χ	1	G	1	1	∞						C C s_X
15	2	0	G	2		2	2					C C G C C G C A s_2
16	2	0	G	3		2	2					C C G G C A s_2
17	1	0	A	1								C G A G A C G A T s_1
18	1	0	A	3								C G A T s_1
19	χ	1	G	2	1							C C C s_X
20	2	0	A	4								C C C C C G C C G G C A s_2
21	2	0	C	5								C G C C C G G C A s_2
22	χ	1	A	3	3							C G C C G C s_X
23	2	0	C	2								C G G C A s_2
24	1	0	A	0								G A C G A T s_1
25	1	0	C	2								G A G A C G A T s_1
26	1	0	C	2								G A T s_1
27	2	0	G	1								G C A s_2
28	χ	1	C	2	0							G C C C s_X
29	2	0	C	3								G C C C G C C G G C A s_2
30	2	0	C	4								G C C C G C C A s_2
31	χ	1	C	2	2							G C C G C s_X
32	2	0	C	1								G C C A s_2
33	1	0	A	0								T s_1
				-1	-1							

Example

$$s_X = \text{ACGCGCC}\$X,$$

$$s_1 = \text{ACGAGACGAT}\$1,$$

$$s_2 = \text{AACGCCGCCGGCA}\$2$$

#	id	S	ebwt	lcp	lcp _X	α	UcLCP			LcLCP	cLCP	Sorted suffixes of S
							χ	1	2			
1	χ	1	C	-1	-1							s_X
2	1	0	T	0								s_1
3	2	0	A	0								s_2
4	2	0	C	0								A s_2
5	2	0	s_2	1								A A C G C G C C G G C A s_2
6	1	0	s_1	1								A C G A T s_1
7	1	0	G	4								A C G A T s_1
8	2	0	A	3								A C G C C G C G G C A s_2
9	χ	1	s_0	4	0							A C G C C G C s_X
10	1	0	G	1								A G A C G A T s_1
11	1	0	G	1								A T s_1
12	χ	1	C	0	0							C s_X
13	2	0	G	1								C A s_2
14	χ	1	G	1	1	∞						C C s_X
15	2	0	G	2		2	2					C C G C C G C A s_2
16	2	0	G	3		2	2					C C G C A s_2
17	1	0	A	1		1	1					C G A G A C G A T s_1
18	1	0	A	3								C G A T s_1
19	χ	1	G	2	1							C C C s_X
20	2	0	A	4								C C C C C G C C G G C A s_2
21	2	0	C	5								C G C C G G C A s_2
22	χ	1	A	3	3							C G C C G C s_X
23	2	0	C	2								C G C A s_2
24	1	0	A	0								G A C G A T s_1
25	1	0	C	2								G A G A C G A T s_1
26	1	0	C	2								G A T s_1
27	2	0	G	1								G C A s_2
28	χ	1	C	2	0							G C C C s_X
29	2	0	C	3								G C C C G C C G G C A s_2
30	2	0	C	4								G C C C G C A s_2
31	χ	1	C	2								G C C G C s_X
32	2	0	C	1								G C C A s_2
33	1	0	A	0								G G C A s_2
				-1	-1							T s_1

Example

$$s_X = ACGCGCC\$_X,$$

$$s_1 = ACGAGACGAT\$_1,$$

$$s_2 = AACGCCGCCGGCA\$_2$$

#	id	S	ebwt	lcp	lcp _X	α	UcLCP			LcLCP	cLCP	Sorted suffixes of S
							χ	1	2			
1	χ	1	C	-1	-1							s_X
2	1	0	T	0								s_1
3	2	0	A	0								s_2
4	2	0	C	0								A s_2
5	2	0	s_2	1								A A C G C G C C G G C A s_2
6	1	0	s_1	1								A C G A T s_1
7	1	0	G	4								A C G A T s_1
8	2	0	A	3								A C G C C G C G G C A s_2
9	χ	1	s_0	4	0							A C G C C G C s_X
10	1	0	G	1								A G A C G A T s_1
11	1	0	G	1								A T s_1
12	χ	1	C	0	0							C s_X
13	2	0	G	1								C A s_2
14	χ	1	G	1	1	∞						C C s_X
15	2	0	G	2		2	2					C C G C C G C A s_2
16	2	0	G	3		2	2					C C G C A s_2
17	1	0	A	1		1	1					C G A G A C G A T s_1
18	1	0	A	3		1	1					C G A T s_1
19	χ	1	G	2	1							C G C C s_X
20	2	0	A	4								C G C C G C G C A s_2
21	2	0	C	5								C G C C G C A s_2
22	χ	1	A	3	3							C G C G C s_X
23	2	0	C	2								C G C A s_2
24	1	0	A	0								G A C G A T s_1
25	1	0	C	2								G A G A C G A T s_1
26	1	0	C	2								G A T s_1
27	2	0	G	1								G C A s_2
28	χ	1	C	2	0							G C C s_X
29	2	0	C	3								G C C G C G C A s_2
30	2	0	C	4								G C C C G C A s_2
31	χ	1	C	2								G C C G C s_X
32	2	0	C	1								G C C A s_2
33	1	0	A	0								T s_1
				-1	-1							

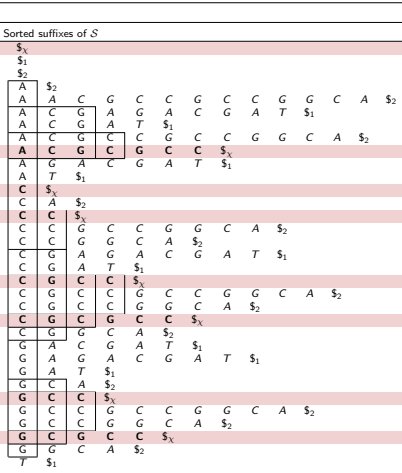
Example

$$s_X = \text{ACGCGCC}\$X,$$

$$s_1 = \text{ACGAGACGAT}\$1,$$

$$s_2 = \text{AACGCCGCCGGCA}\$2$$

#	id	S	ebwt	lcp	lcp _X	α	UcLCP			LcLCP	cLCP	Sorted suffixes of S
							X	1	2			
1	X	1	C	-1	-1	∞						\$ _X
2	1	0	T	0		0	0					\$ ₁
3	2	0	A	0		0	0					\$ ₂
4	2	0	C	0		0	0					
5	2	0	\$ ₂	1		0	0					
6	1	0	\$ ₁	1		0	0					
7	1	0	G	4		0	0					
8	2	0	A	3		0	0					
9	X	1	\$ ₀	4	0	∞						
10	1	0	G	1		1	1					
11	1	0	G	1		1	1					
12	X	1	C	0	0	∞						
13	2	0	G	1		1	1					
14	X	1	G	1	1	∞						
15	2	0	G	2		2	2					
16	2	0	G	3		2	2					
17	1	0	A	1		1	1					
18	1	0	A	3		1	1					
19	X	1	G	2	1	∞						
20	2	0	A	4		4	4					
21	2	0	C	5		4	4					
22	X	1	A	3	3	∞						
23	2	0	C	2		2	2					
24	1	0	A	0		0	0					
25	1	0	C	2		0	0					
26	1	0	C	2		0	0					
27	2	0	G	1		0	0					
28	X	1	C	2	0	∞						
29	2	0	C	3		3	3					
30	2	0	C	4		3	3					
31	X	1	C	2	2	∞						
32	2	0	C	1		1	1					
33	1	0	A	0		0	0					



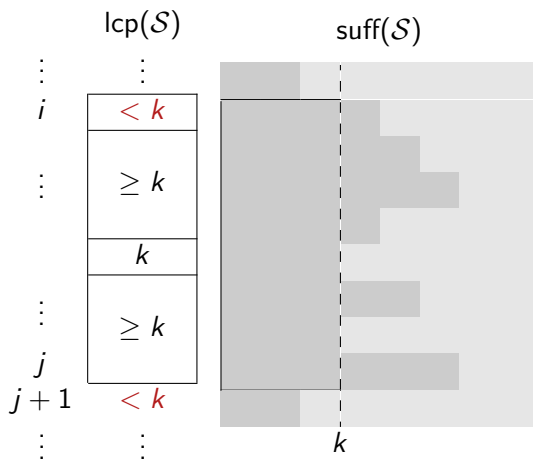
LcLCP[i_r][χ] Computation

i	$\text{lcp}(S)$	$\text{id}(S)$
\vdots	\vdots	\vdots
$\chi_1 = \text{prev}(i_r, \chi)$		χ
\vdots	\vdots	\top
		$\neq \chi$
		\perp
i_r		r
\vdots	\vdots	\top
		$\neq \chi$
		\perp
$\chi_2 = \text{next}(i_r, \chi)$		χ
\vdots	\vdots	\vdots

Problem: LcLCP[i_r][χ] computation would require to look **forward** and to store many intermediate values.

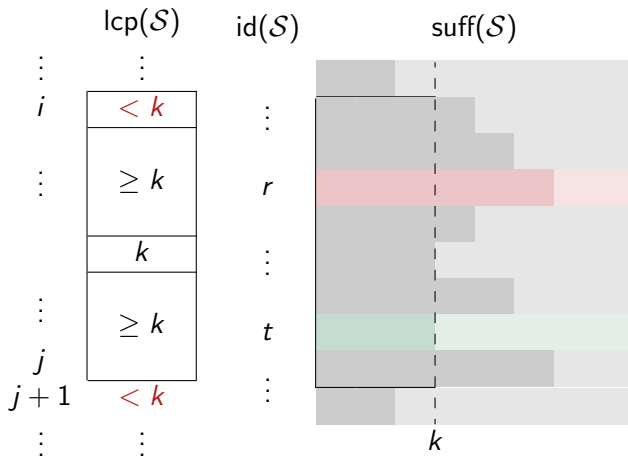
Colored k -lcp interval

A k -lcp interval is an interval $[i, j]$ such that:



Colored k -lcp interval

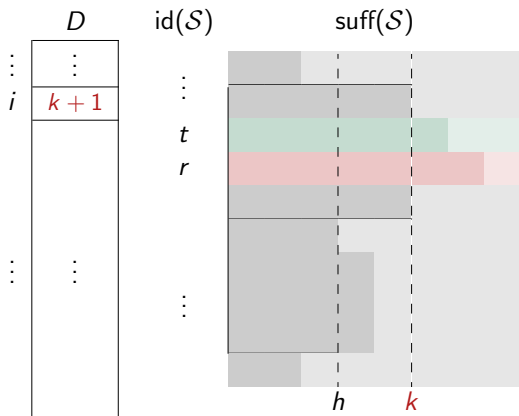
A colored k -lcp interval is an interval $[i, j]$ such that:



s_r and s_t belong one to S^0 , the other to S^1

Array D

Denote with D the $(N + 1)$ -integer array such that $D[i] = k + 1$ if a colored k -lcp interval starts at position i and for every colored h -lcp interval starting at position i then $h \leq k - 1$.



LcLCP[i_r][χ] Computation

Theorem

For any $1 \leq i_r \leq N$ such that $\text{id}(\mathcal{S})[i_r] = r$,

- if $\text{LCP}(\chi_1, i_r) > \text{LCP}(\chi_1, \chi_2)$ then $\text{LCP}(i_r, \chi_2) = \text{LCP}(\chi_1, \chi_2)$
- otherwise

$$\text{LCP}(i_r, \chi_2) = \max\{\max\{D[x] : \chi_1 < x \leq i_r\} - 1, \text{LCP}(\chi_1, \chi_2)\}$$

LcLCP[i_r][χ] Computation

Theorem

For any $1 \leq i_r \leq N$ such that $\text{id}(S)[i_r] = r$,

- if $\text{LCP}(\chi_1, i_r) > \text{LCP}(\chi_1, \chi_2)$ then $\text{LCP}(i_r, \chi_2) = \text{LCP}(\chi_1, \chi_2)$
- otherwise

$$\text{LCP}(i_r, \chi_2) = \max\{\max\{D[x] : \chi_1 < x \leq i_r\} - 1, \text{LCP}(\chi_1, \chi_2)\}$$

▷ Keep track of the maximum D value since the beginning of each χ -interval.

	i	$\text{lcp}(S)$	D	$\text{id}(S)$
	\vdots	\vdots	\vdots	\vdots
$\zeta \leftarrow 0$	$\chi_1 = \text{prev}(i_r, \chi)$			χ
	\vdots	\vdots	\vdots	\top
	i_r			$\neq \chi$
	\vdots	\vdots	\vdots	\perp
	$\chi_2 = \text{next}(i_r, \chi)$			r
	\vdots	\vdots	\vdots	\top
				$\neq \chi$
				\perp
				χ
	\vdots	\vdots	\vdots	\vdots

LcLCP[i_r][χ] Computation

Theorem

For any $1 \leq i_r \leq N$ such that $\text{id}(S)[i_r] = r$,

- if $\text{LCP}(\chi_1, i_r) > \text{LCP}(\chi_1, \chi_2)$ then $\text{LCP}(i_r, \chi_2) = \text{LCP}(\chi_1, \chi_2)$
- otherwise

$$\text{LCP}(i_r, \chi_2) = \max\{\max\{D[x] : \chi_1 < x \leq i_r\} - 1, \text{LCP}(\chi_1, \chi_2)\}$$

▷ Keep track of the maximum D value since the beginning of each χ -interval.

	i	$\text{lcp}(S)$	D	$\text{id}(S)$
	\vdots	\vdots	\vdots	\vdots
$\chi_1 = \text{prev}(i_r, \chi)$				χ
	\vdots	\vdots	\vdots	\top
				$\neq \chi$
				\perp
$\zeta \leftarrow \max D(\chi_1, i_r)$	i_r			r
	\vdots	\vdots	\vdots	\top
				$\neq \chi$
				\perp
$\chi_2 = \text{next}(i_r, \chi)$				χ
	\vdots	\vdots	\vdots	\vdots

Example

#	id	S	ebwt	lcp	D	lcp _x	α	ζ	UcLCP			LcLCP			cLCP	Sorted suffixes of S
									χ	1	2	χ	1	2		
1	χ	1	C	-1	0	-1	∞									\$ _x
2	1	0	T	0	0		0	0								\$ ₁
3	2	0	A	0	0		0	0								\$ ₂
4	2	0	C	0	2		0	0								A \$ ₂
5	2	0	\$ ₂	1	0		0	0								A A C G C C G C C G G C A \$ ₂
6	1	0	\$ ₁	1	4		0	0								A C G C A G A C G A T \$ ₁
7	1	0	G	4	0		0	0								A C G A T \$ ₁
8	2	0	A	3	5		0	0								A C G C C G C C C G G C A \$ ₂
9	χ	1	\$ ₀	4	0	0	∞									A C G C G C C \$ _x
10	1	0	G	1	0		1	1								A G A C G A T \$ ₁
11	1	0	G	1	0		1	1								A T \$ ₁
12	χ	1	C	0	2	0	∞									C \$ _x
13	2	0	G	1	0		1	1								C A \$ ₂
14	χ	1	G	1	3	1	∞									C C \$ _x
15	2	0	G	2	0		2	2								C C G C C G G C C A \$ ₂
16	2	0	G	3	0		2	2								C C G G C A A \$ ₂
17	1	0	A	1	3		1	1								C G A G A C G A T \$ ₁
18	1	0	A	3	0		1	1								C G A T \$ ₁
19	χ	1	G	2	5	1	∞									C G C C \$ _x
20	2	0	A	4	0		4	4								C G C C G C C G G C A \$ ₂
21	2	0	C	5	0		4	4								C G C C G G C A \$ ₂
22	χ	1	A	3	0	3	∞									C G C G C C \$ _x
23	2	0	C	2	0		2	2								C G G C A \$ ₂
24	1	0	A	0	2		0	0								G A C G A T \$ ₁
25	1	0	C	2	0		0	0								G A G A C G A T \$ ₁
26	1	0	C	2	0		0	0								G A T \$ ₁
27	2	0	G	1	3		α	0								G C A \$ ₂
28	χ	1	C	2	4	0	∞									G C C \$ _x
29	2	0	C	3	0		3	3								G C C G C C G G C A \$ ₂
30	2	0	C	4	0		3	3								G C C G G C A \$ ₂
31	χ	1	C	2	0	2	∞									G C G C C \$ _x
32	2	0	C	1	0		1	1								G G C A \$ ₂
33	1	0	A	0	0		0	0								T \$ ₁

Example

#	id	S	ebwt	lcp	D	lcp _x	α	ζ	UcLCP			LcLCP			cLCP	Sorted suffixes of S
									χ	1	2	χ	1	2		
1	χ	1	C	-1	0	-1	∞									\$ _x
2	1	0	T	0	0		0	0								\$ ₁
3	2	0	A	0	0		0	0								\$ ₂
4	2	0	C	0	2		0	0								A \$ ₂
5	2	0	\$ ₂	1	0		0	0								A C \$ ₂
6	1	0	\$ ₁	1	4		0	0								A C G A A G A C G A T \$ ₁ C A \$ ₂
7	1	0	G	4	0		0	0								A C G A T \$ ₁
8	2	0	A	3	5		0	0								A C G C C G C C C G G C A \$ ₂
9	χ	1	\$ ₀	4	0	0	∞									A C G C G C C \$ _x
10	1	0	G	1	0		1	1								A G A C G A T \$ ₁
11	1	0	G	1	0		1	1								A T \$ ₁
12	χ	1	C	0	2	0	∞									C \$ _x
13	2	0	G	1	0		1	1								C A \$ ₂
14	χ	1	G	1	3	1	∞	0								C C \$ _x
15	2	0	G	2	0		2	2								C C G C C G G C C A \$ ₂
16	2	0	G	3	0		2	2								C C G G C A A \$ ₂
17	1	0	A	1	3		1	1								C G A G A C G A T \$ ₁
18	1	0	A	3	0		1	1								C G A T \$ ₁
19	χ	1	G	2	5	1	∞									C G C C \$ _x
20	2	0	A	4	0		4	4								C G C C G C C G G C A \$ ₂
21	2	0	C	5	0		4	4								C G C C G G C A \$ ₂
22	χ	1	A	3	0	3	∞									C G C G C C \$ _x
23	2	0	C	2	0		2	2								C G G C A \$ ₂
24	1	0	A	0	2		0	0								G A C G A T \$ ₁
25	1	0	C	2	0		0	0								G A G A C G A T \$ ₁
26	1	0	C	2	0		0	0								G A T \$ ₁
27	2	0	G	1	3		α	0								G C A \$ ₂
28	χ	1	C	2	4	0	∞									G C C \$ _x
29	2	0	C	3	0		3	3								G C C G C C G G C A \$ ₂
30	2	0	C	4	0		3	3								G C C G G C A \$ ₂
31	χ	1	C	2	0	2	∞									G C G C C \$ _x
32	2	0	C	1	0		1	1								G G C A \$ ₂
33	1	0	A	0	0		0	0								T \$ ₁

Example

#	id	S	ebwt	lcp	D	lcp _x	α	ζ	UcLCP			LcLCP			cLCP	Sorted suffixes of S
									χ	1	2	χ	1	2		
1	χ	1	C	-1	0	-1	∞									\$ _x
2	1	0	T	0	0		0	0								\$ ₁
3	2	0	A	0	0		0	0								\$ ₂
4	2	0	C	0	2		0	0								A \$ ₂
5	2	0	\$ ₂	1	0		0	0								A A C G C C G C C G G C A \$ ₂
6	1	0	\$ ₁	1	4		0	0								A C G C A G A C G A T \$ ₁
7	1	0	G	4	0		0	0								A C G A T \$ ₁
8	2	0	A	3	5		0	0								A C G C C G C C C G G C A \$ ₂
9	χ	1	\$ ₀	4	0	0	∞									A C G C C G C C C \$ _x
10	1	0	G	1	0		1	1								A G A C G A T \$ ₁
11	1	0	G	1	0		1	1								A T \$ ₁
12	χ	1	C	0	2	0	∞									C \$ _x
13	2	0	G	1	0		1	1								C A \$ ₂
14	χ	1	G	1	3	1	∞	0								C C \$ _x
15	2	0	G	2	0		2	2		1						C C G C C G G C C A \$ ₂
16	2	0	G	3	0		2	2								C C G G C A A \$ ₂
17	1	0	A	1	3		1	1								C G A G A C G A T \$ ₁
18	1	0	A	3	0		1	1								C G A T \$ ₁
19	χ	1	G	2	5	1	∞									C G C C C \$ _x
20	2	0	A	4	0		4	4								C G C C G C C G C C G G C A \$ ₂
21	2	0	C	5	0		4	4								C G C C C G G C C A \$ ₂
22	χ	1	A	3	0	3	∞									C G C C G C C \$ _x
23	2	0	C	2	0		2	2								C G G C A \$ ₂
24	1	0	A	0	2		0	0								G A C G A T \$ ₁
25	1	0	C	2	0		0	0								G A G A C G A T \$ ₁
26	1	0	C	2	0		0	0								G A T \$ ₁
27	2	0	G	1	3		α	0								G C A \$ ₂
28	χ	1	C	2	4	0	∞									G C C C \$ _x
29	2	0	C	3	0		3	3								G C C G C C G G C A \$ ₂
30	2	0	C	4	0		3	3								G C C G G C A \$ ₂
31	χ	1	C	2	0	2	∞									G C G C C \$ _x
32	2	0	C	1	0		1	1								G G C A \$ ₂
33	1	0	A	0	0		0	0								T \$ ₁

Example

#	id	S	ebwt	lcp	D	lcp _x	α	ζ	UcLCP			LcLCP			cLCP	Sorted suffixes of S
									χ	1	2	χ	1	2		
1	χ	1	C	-1	0	-1	∞									\$ _x
2	1	0	T	0	0		0	0								\$ ₁
3	2	0	A	0	0		0	0								\$ ₂
4	2	0	C	0	2		0	0								A \$ ₂
5	2	0	\$ ₂	1	0		0	0								A A C G C C G C C G G C A \$ ₂
6	1	0	\$ ₁	1	4		0	0								A C G C A G A C G A T \$ ₁
7	1	0	G	4	0		0	0								A C G A T \$ ₁
8	2	0	A	3	5		0	0								A C G C C G C C C G G C A \$ ₂
9	χ	1	\$ ₀	4	0	0	∞									A C G C C G C C C \$ _x
10	1	0	G	1	0		1	1								A G A C G A T \$ ₁
11	1	0	G	1	0		1	1								A T \$ ₁
12	χ	1	C	0	2	0	∞									C \$ _x
13	2	0	G	1	0		1	1								C A \$ ₂
14	χ	1	G	1	3	1	∞	0								C C \$ _x
15	2	0	G	2	0		2	0	2			1				C C G C C G G C C G G C A \$ ₂
16	2	0	G	3	0		2	0	2			1				C C G G C A A \$ ₂
17	1	0	A	1	3		1	1	1							C G A G A C G A T \$ ₁
18	1	0	A	3	0		1	1	1							C G A T \$ ₁
19	χ	1	G	2	5	1	∞									C G C C C \$ _x
20	2	0	A	4	0		4	4	4							C G C C C G C C G G C A \$ ₂
21	2	0	C	5	0		4	4	4							C G C C C G G C A \$ ₂
22	χ	1	A	3	0	3	∞									C G C C G C C \$ _x
23	2	0	C	2	0		2	2	2							C G G C A \$ ₂
24	1	0	A	0	2		0	0	0							G A C G A T \$ ₁
25	1	0	C	2	0		0	0	0							G A G A C G A T \$ ₁
26	1	0	C	2	0		0	0	0							G A T \$ ₁
27	2	0	G	1	3		0	0	0							G C A \$ ₂
28	χ	1	C	2	4	0	∞									G C C C \$ _x
29	2	0	C	3	0		3	3	3							G C C G C C G G C A \$ ₂
30	2	0	C	4	0		3	3	3							G C C G G C A \$ ₂
31	χ	1	C	2	0	2	∞									G C G C C \$ _x
32	2	0	C	1	0		1	1	1							G G C A \$ ₂
33	1	0	A	0	0		0	0	0							T \$ ₁

Example

#	id	S	ebwt	lcp	D	lcp _x	α	ζ	UcLCP			LcLCP			cLCP	Sorted suffixes of S
									χ	1	2	χ	1	2		
1	χ	1	C	-1	0	-1	∞									\$ _x
2	1	0	T	0	0		0	0								\$ ₁
3	2	0	A	0	0		0	0								\$ ₂
4	2	0	C	0	2		0	0								A \$ ₂
5	2	0	\$ ₂	1	0		0	0								A A C G C C G C C G G C A \$ ₂
6	1	0	\$ ₁	1	4		0	0								A C G C A G A C G A T \$ ₁
7	1	0	G	4	0		0	0								A C G A T \$ ₁
8	2	0	A	3	5		0	0								A C G C C G C C C G G C A \$ ₂
9	χ	1	\$ ₀	4	0	0	∞									A C G C G C C \$ _x
10	1	0	G	1	0		1	1								A G A C G A T \$ ₁
11	1	0	G	1	0		1	1								A T \$ ₁
12	χ	1	C	0	2	0	∞									C \$ _x
13	2	0	G	1	0		1	1								C A \$ ₂
14	χ	1	G	1	3	1	∞	0								C C \$ _x
15	2	0	G	2	0		2	0	2			1				C C G C C G G C C G G C A \$ ₂
16	2	0	G	3	0		2	0	2			1				C C G G C A A \$ ₂
17	1	0	A	1	3		1	2	1			2				C C G A G A C G A T \$ ₁
18	1	0	A	3	0		1	1	1			2				C G A T \$ ₁
19	χ	1	G	2	5	1	∞									C G C C \$ _x
20	2	0	A	4	0		4	4	4							C G C C G C C G C C G G C A \$ ₂
21	2	0	C	5	0		4	4	4							C G C C G G C A \$ ₂
22	χ	1	A	3	0	3	∞									C G C G C C \$ _x
23	2	0	C	2	0		2	2	2							C G G C A \$ ₂
24	1	0	A	0	2		0	0	0							G A C G A T \$ ₁
25	1	0	C	2	0		0	0	0							G A G A C G A T \$ ₁
26	1	0	C	2	0		0	0	0							G A T \$ ₁
27	2	0	G	1	3		α	α	0							G C A \$ ₂
28	χ	1	C	2	4	0	∞									G C C \$ _x
29	2	0	C	3	0		3	3	3							G C C G C C G G C A \$ ₂
30	2	0	C	4	0		3	3	3							G C C G G C A \$ ₂
31	χ	1	C	2	0	2	∞									G C G C C \$ _x
32	2	0	C	1	0		1	1	1							G G C A \$ ₂
33	1	0	A	0	0		0	0	0							T \$ ₁

Example

#	id	S	ebwt	lcp	D	lcp _x	α	ζ	UcLCP			LcLCP			cLCP	Sorted suffixes of S
									χ	1	2	χ	1	2		
1	χ	1	C	-1	0	-1	∞									\$ _x
2	1	0	T	0	0		0	0								\$ ₁
3	2	0	A	0	0		0	0								\$ ₂
4	2	0	C	0	2		0	0								A \$ ₂
5	2	0	\$ ₂	1	0		0	0								A A C G C C G C C G G C A \$ ₂
6	1	0	\$ ₁	1	4		0	0								A C G A A G A C G A T \$ ₁
7	1	0	G	4	0		0	0								A C G A T \$ ₁
8	2	0	A	3	5		0	0								A C G C C G C C C G G C A \$ ₂
9	χ	1	\$ ₀	4	0	0	∞									A C G C G C C \$ _x
10	1	0	G	1	0		1	1								A G A C G A T \$ ₁
11	1	0	G	1	0		1	1								A T \$ ₁
12	χ	1	C	0	2	0	∞									C \$ _x
13	2	0	G	1	0		1	1								C A \$ ₂
14	χ	1	G	1	3	1	∞	0								C C \$ _x
15	2	0	G	2	0		2	0	2			1				C C G C C G G C C G G C A \$ ₂
16	2	0	G	3	0		2	0	2			1				C C G G C A A \$ ₂
17	1	0	A	1	3		1	2	1			2				C G A G A C G A T \$ ₁
18	1	0	A	3	0		1	2	1			2				C G A T \$ ₁
19	χ	1	G	2	5	1	∞									C G C C \$ _x
20	2	0	A	4	0		4	4								C G C C G C C G C C G G C A \$ ₂
21	2	0	C	5	0		4	4								C G C C G G C C A \$ ₂
22	χ	1	A	3	0	3	∞									C G C G C C \$ _x
23	2	0	C	2	0		2	2								C G G C A \$ ₂
24	1	0	A	0	2		0	0								G A C G A T \$ ₁
25	1	0	C	2	0		0	0								G A G A C G A T \$ ₁
26	1	0	C	2	0		0	0								G A T \$ ₁
27	2	0	G	1	3		α	0								G C A \$ ₂
28	χ	1	C	2	4	0	∞									G C C \$ _x
29	2	0	C	3	0		3	3								G C C G C C G G C A \$ ₂
30	2	0	C	4	0		3	3								G C C G G C A \$ ₂
31	χ	1	C	2	0	2	∞									G C G C C \$ _x
32	2	0	C	1	0		1	1								G G C A \$ ₂
33	1	0	A	0	0		0	0								T \$ ₁

Example

#	id	S	ebwt	lcp	D	lcp _x	α	ζ	UcLCP			LcLCP			cLCP	Sorted suffixes of S
									X	1	2	X	1	2		
1	X	1	C	-1	0	-1	∞	0								S _X
2	1	0	T	0	0		0	0				0				S ₁
3	2	0	A	0	0		0	0				0				S ₂
4	2	0	C	0	2		0	1				0				A
5	2	0	S ₂	1	0		0	1				1				A
6	1	0	S ₁	1	4		0	3				0				A
7	1	0	G	4	0		0	3				0				A
8	2	0	A	3	5		0	4				0				A
9	X	1	S ₀	4	0	0	∞	0								A
10	1	0	G	1	0		1	0	1			0				A
11	1	0	G	1	0		1	0	1			0				A
12	X	1	C	0	2	0	∞	0								A
13	2	0	G	1	0		1	0	1			1				A
14	X	1	G	1	3	1	∞	0								A
15	2	0	G	2	0		2	0	2			1				A
16	2	0	G	3	0		2	0	2			1				A
17	1	0	A	1	3		1	2	1			2				A
18	1	0	A	3	0		1	2	1			2				A
19	X	1	G	2	5	1	∞	0								A
20	2	0	A	4	0		4	0	4			3				A
21	2	0	C	5	0		4	0	4			3				A
22	X	1	A	3	0	3	∞	0								A
23	2	0	C	2	0		2	0	2			0				A
24	1	0	A	0	2		0	1	0			1				A
25	1	0	C	2	0		0	1	0			1				A
26	1	0	C	2	0		0	1	0			1				A
27	2	0	G	1	3		0	2	0			2				A
28	X	1	C	2	4	0	∞	0								A
29	2	0	C	3	0		3	2	3			2				A
30	2	0	C	4	0		3	3	3			2				A
31	X	1	C	2	0	2	∞	0								A
32	2	0	C	1	0		1	0	1			0				A
33	1	0	A	0	0		0	0	0			0				A

Missing $\text{LcLCP}[\chi_1][r]$, $\text{UcLCP}[\chi_2][r]$ values

For each χ -interval $[\chi_1, \chi_2]$, we are able to compute $\text{UcLCP}[i_r][\chi]$ and $\text{LcLCP}[i_r][\chi]$, but not $\text{LcLCP}[\chi_1][r]$ and $\text{UcLCP}[\chi_2][r]$ values in a single sequential scan.

i	$\text{lcp}(S)$	$\text{id}(S)$
\vdots	\vdots	\vdots
χ_1		χ
\vdots	\vdots	\top
		$\neq \chi$
		\perp
i_r		r
\vdots	\vdots	\top
		$\neq \chi$
		\perp
χ_2		χ
\vdots	\vdots	\vdots

What is the reason of this asymmetry?

Missing LcLCP[χ_1][r], UcLCP[χ_2][r] values (2)

- The segmentation of $\text{lcp}(\mathcal{S})$ in χ -intervals, induced by s_χ suffixes, is a *shared* reference amongst all the suffixes of any other color r .

i	$\text{lcp}(\mathcal{S})$	$\text{id}(\mathcal{S})$
\vdots	\vdots	\vdots
χ_1		χ
\vdots	\vdots	\top
		$\neq \chi$
		\perp
i_r		r
\vdots	\vdots	\top
		$\neq \chi$
		\perp
χ_2		χ
\vdots	\vdots	\vdots

- However, no $\text{lcp}(\mathcal{S})$ segmentation in r -intervals (induced by $s_r \in S^1$ suffixes) is maintained for suffixes of s_χ during the sequential scan.

LcLCP[χ_1][r] Computation

- In fact, it is easy to compute $\text{LcLCP}[\chi_1][r]$ if there exists a r -colored suffix in (χ_1, χ_2) :

i	$\text{lcp}(S)$	$\text{id}(S)$
\vdots	\vdots	\vdots
χ_1		χ
\vdots	\vdots	\top
$i_r = \text{next}(\chi_1, r)$		$\neq \chi, r$
		\perp
		r
\vdots	\vdots	\top
		$\neq \chi$
		\perp
χ_2		χ
\vdots	\vdots	\vdots

$$\text{LcLCP}[\chi_1][r] = \text{UcLCP}[i_r][\chi]$$

LcLCP[χ_1][r] Computation - Backward Propagation

- However, what if such a suffix does not exist?

i	$\text{lcp}(S)$	$\text{id}(S)$
\vdots	\vdots	\vdots
χ_1		χ
\vdots	\vdots	\top
		$\neq \chi, r$
		\perp
χ_2		χ
\vdots	\vdots	\top
		$\neq r$
		\perp
$i_r = \text{next}(\chi_2, r)$		r
\vdots	\vdots	\vdots

We need to look beyond χ_2 and **backward propagate** $\text{LcLCP}[\chi_2][r]$:

$$\text{LcLCP}[\chi_1][r] = \min\{\text{LCP}(\chi_1, \chi_2), \text{UcLCP}[i_r][\chi]\}$$

Theorem

Let \mathcal{S} a collection of m strings partitioned into \mathcal{S}^0 and \mathcal{S}^1 , let $s_\chi \in \mathcal{S}^0$. Given $\text{id}(\mathcal{S})$, $\text{lcp}(\mathcal{S})$ and $\text{lcp}(s_\chi)$, $\text{cLCP}(\mathcal{S})$ can be computed by sequential scans in $\mathcal{O}(N + m|s_\chi|)$ time and $\mathcal{O}(m + \max \text{lcp}(\mathcal{S}))$ space.

- Array D : in $\Theta(N)$ time and $\mathcal{O}(\max \text{lcp}(\mathcal{S}))$ space.
- $\text{UcLCP}[i_r][\chi]/\text{LcLCP}[i_r][\chi]$: sequentially (forward) in $\Theta(N)$ time and $\mathcal{O}(1)$ space.
- $\text{UcLCP}[\chi_2][r]/\text{LcLCP}[\chi_1][r]$: add $\mathcal{O}(m|s_\chi|)$ time and $\mathcal{O}(m)$ space for forward and backward propagation.

Theorem

Let \mathcal{S} a collection of m strings partitioned into \mathcal{S}^0 and \mathcal{S}^1 , let $s_\chi \in \mathcal{S}^0$. Given $\text{id}(\mathcal{S})$, $\text{lcp}(\mathcal{S})$ and $\text{lcp}(s_\chi)$, $\text{cLCP}(\mathcal{S})$ can be computed by sequential scans in $\mathcal{O}(N + m|s_\chi|)$ time and $\mathcal{O}(m + \max \text{lcp}(\mathcal{S}))$ space.

- Array D: in $\Theta(N)$ time and $\mathcal{O}(\max \text{lcp}(\mathcal{S}))$ space.
- $\text{UcLCP}[i_r][\chi]/\text{LcLCP}[i_r][\chi]$: sequentially (forward) in $\Theta(N)$ time and $\mathcal{O}(1)$ space.
- $\text{UcLCP}[\chi_2][r]/\text{LcLCP}[\chi_1][r]$: add $\mathcal{O}(m|s_\chi|)$ time and $\mathcal{O}(m)$ space for forward and backward propagation.

...affordable for solving the (1-vs-all) multi-string ACS problem, BUT unsuitable for an all-vs-all scenario!

Can we compute $UcLCP[\chi_2][r]/LcLCP[\chi_1][r]$ without relying on backward-forward propagation?

Avoiding backward-forward propagation

- Idea:** Get rid of the $\text{lcp}(\mathcal{S})$ χ -induced segmentation as a shared reference for all the other suffixes, and keep track of *all* possible t -intervals, for $t \in [1, m]$.

i	$\text{lcp}(\mathcal{S})$	$\text{id}(\mathcal{S})$
\vdots	\vdots	\vdots
$t_1 = \text{prev}(i_r, t)$	\vdots	t
\vdots	\vdots	\top
\vdots	\vdots	$\neq t$
i_r	\vdots	\perp
\vdots	\vdots	r
\vdots	\vdots	\top
\vdots	\vdots	$\neq t$
\vdots	\vdots	\perp
$t_2 = \text{next}(i_r, t)$	\vdots	t
\vdots	\vdots	\vdots

Avoiding backward-forward propagation (2)

- Maintain $\alpha[1, m]$ such that $\alpha[t]$ keeps track of the minimum lcp value since the beginning of each t -interval:

$$\alpha[t] = \min\{\text{lcp}(\mathcal{S})[x] : x \in (\text{prev}(i_r, t), i_r]\}$$

- Maintain $\zeta[1, m]$ such that $\zeta[t]$ keeps track of the maximum D_t value since the beginning of each t -interval:

$$\zeta[t] = \max\{D_t[x] : x \in (\text{prev}(i_r, t), i_r]\}$$

Avoiding backward-forward propagation (3)

- For the computation of $\text{UcLCP}[i_r][t]$, we have:

$$\text{UcLCP}[i_r][t] = \min\{\text{lcp}(\mathcal{S})[x] : x \in (t_1, i_r)\} = \alpha[t]$$

- For the computation of $\text{LcLCP}[i_r][t]$, we use Theorem 1:

Theorem

For any $1 \leq i_r \leq N$ such that $\text{id}(\mathcal{S})[i_r] = r$,

- if $\text{LCP}(t_1, i_r) > \text{LCP}(t_1, t_2)$ then $\text{LCP}(i_r, t_2) = \text{LCP}(t_1, t_2)$
- otherwise

$$\text{LCP}(i_r, t_2) = \max\{\max\{\text{D}_t[x] : t_1 < x \leq i_r\} - 1, \text{LCP}(t_1, t_2)\}$$

Avoiding backward-forward propagation (3)

- For the computation of $\text{UcLCP}[i_r][t]$, we have:

$$\text{UcLCP}[i_r][t] = \min\{\text{lcp}(\mathcal{S})[x] : x \in (t_1, i_r)\} = \alpha[t]$$

- For the computation of $\text{LcLCP}[i_r][t]$, we use Theorem 1:

Theorem

For any $1 \leq i_r \leq N$ such that $\text{id}(\mathcal{S})[i_r] = r$,

- if $\alpha[t] > \text{LCP}(t_1, t_2)$ then $\text{LCP}(i_r, t_2) = \text{LCP}(t_1, t_2)$
- otherwise

$$\text{LCP}(i_r, t_2) = \max\{\zeta[t] - 1, \text{LCP}(t_1, t_2)\}$$

No more forward-backward propagation: purely sequential $\mathcal{O}(Nm)$ time solution for an all-vs-all ACS computation scenario!

Matching Statistics

Given two strings s_0 , s_1 , the **matching statistics** between of s_0 vs. s_1 is a $|s_0|$ -integer array such that, for any position i of s_0 , stores the length of the longest prefix of the suffix of s_0 starting at position i that is also substring of s_1 .

s_0	A	C	G	C	G	C	C				
s_1	A	C	G	A	G	A	C	G	A	T	

$MS(s_0, s_1)$

Matching Statistics

Given two strings s_0 , s_1 , the **matching statistics** between of s_0 vs. s_1 is a $|s_0|$ -integer array such that, for any position i of s_0 , stores the length of the longest prefix of the suffix of s_0 starting at position i that is also substring of s_1 .

	↓										
s_0	A	C	G	C	G	C	C				
s_1	A	C	G	A	G	A	C	G	A	T	

$MS(s_0, s_1)$ 3

Matching Statistics

Given two strings s_0 , s_1 , the **matching statistics** between of s_0 vs. s_1 is a $|s_0|$ -integer array such that, for any position i of s_0 , stores the length of the longest prefix of the suffix of s_0 starting at position i that is also substring of s_1 .

		↓									
s_0	A	C	G	C	G	C	C				
s_1	A	C	G	A	G	A	C	G	A	T	

$MS(s_0, s_1)$	3	2								
----------------	---	---	--	--	--	--	--	--	--	--

Matching Statistics

Given two strings s_0 , s_1 , the **matching statistics** between s_0 vs. s_1 is a $|s_0|$ -integer array such that, for any position i of s_0 , stores the length of the longest prefix of the suffix of s_0 starting at position i that is also substring of s_1 .

			↓								
s_0	A	C	G	C	G	C	C				
s_1	A	C	G	A	G	A	C	G	A	T	
$MS(s_0, s_1)$	3	2	1								

Matching Statistics

Given two strings s_0 , s_1 , the **matching statistics** between of s_0 vs. s_1 is a $|s_0|$ -integer array such that, for any position i of s_0 , stores the length of the longest prefix of the suffix of s_0 starting at position i that is also substring of s_1 .

				↓							
s_0	A	C	G	C	G	C	C				
s_1	A	C	G	A	G	A	C	G	A	T	
$MS(s_0, s_1)$	3	2	1	2							

Matching Statistics

Given two strings s_0 , s_1 , the **matching statistics** between of s_0 vs. s_1 is a $|s_0|$ -integer array such that, for any position i of s_0 , stores the length of the longest prefix of the suffix of s_0 starting at position i that is also substring of s_1 .

					↓						
s_0	A	C	G	C	G	C	C				
s_1	A	C	G	A	G	A	C	G	A	T	
$MS(s_0, s_1)$	3	2	1	2	1						

Matching Statistics

Given two strings s_0 , s_1 , the **matching statistics** between of s_0 vs. s_1 is a $|s_0|$ -integer array such that, for any position i of s_0 , stores the length of the longest prefix of the suffix of s_0 starting at position i that is also substring of s_1 .

						↓					
s_0	A	C	G	C	G	C	C				
s_1	A	C	G	A	G	A	C	G	A	T	
$MS(s_0, s_1)$	3	2	1	2	1	1					

Matching Statistics

Given two strings s_0 , s_1 , the **matching statistics** between of s_0 vs. s_1 is a $|s_0|$ -integer array such that, for any position i of s_0 , stores the length of the longest prefix of the suffix of s_0 starting at position i that is also substring of s_1 .

							↓			
s_0	A	C	G	C	G	C	C			
s_1	A	C	G	A	G	A	C	G	A	T
$MS(s_0, s_1)$	3	2	1	2	1	1	1			

Average Common Substring (ACS) measure

Given two strings s_r and s_t over the alphabet Σ of size σ , ACS is computed by proceeding in the following steps:

- 1 $MS(s_r, s_t)$

Average Common Substring (ACS) measure

Given two strings s_r and s_t over the alphabet Σ of size σ , ACS is computed by proceeding in the following steps:

① $MS(s_r, s_t)$

②
$$\text{Score}(s_r, s_t) = \frac{\sum_{j=1}^{|s_r|} MS(s_r, s_t)[j]}{|s_r|}$$

Average Common Substring (ACS) measure

Given two strings s_r and s_t over the alphabet Σ of size σ , ACS is computed by proceeding in the following steps:

① $MS(s_r, s_t)$

② $Score(s_r, s_t) = \frac{\sum_{j=1}^{|s_r|} MS(s_r, s_t)[j]}{|s_r|}$

③ $Norm(Score(s_r, s_t)) = \frac{\log_{\sigma} |s_t|}{Score(s_r, s_t)} - \frac{2 \log_{\sigma} |s_r|}{|s_r| + 1}$

$$ACS(s_r, s_t) = \frac{Norm(Score(s_r, s_t)) + Norm(Score(s_t, s_r))}{2}$$

Average Common Substring (ACS) measure

Given two strings s_r and s_t over the alphabet Σ of size σ , ACS is computed by proceeding in the following steps:

① $MS(s_r, s_t)$

② $Score(s_r, s_t) = \frac{\sum_{j=1}^{|s_r|} MS(s_r, s_t)[j]}{|s_r|}$

③ $Norm(Score(s_r, s_t)) = \frac{\log_{\sigma} |s_t|}{Score(s_r, s_t)} - \frac{2 \log_{\sigma} |s_r|}{|s_r| + 1}$

$$ACS(s_r, s_t) = \frac{Norm(Score(s_r, s_t)) + Norm(Score(s_t, s_r))}{2}$$

Multi-String ACS Problem

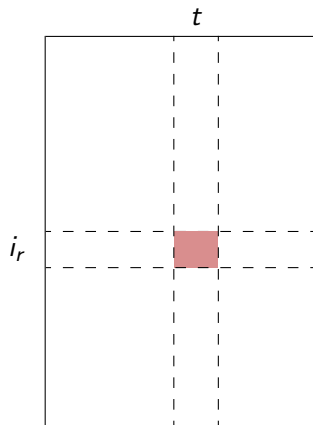
1-vs-all: Compute the pairwise ACS measure between a given string $s_x \in \mathcal{S}^0$ and each string of a set \mathcal{S}^1 of m strings, simultaneously.

all-vs-all: Compute the pairwise ACS measure for any pair of strings $s_r, s_t \in \mathcal{S}$.

MS with cLCP (1)

Proposition

$MS(s_r, s_t)$ is a permutation of the values in $cLCP(\mathcal{S})$ related to suffixes of s_r versus s_t .

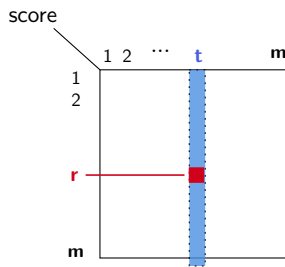


j_r
 $MS(s_r, s_t)$

j_r the initial position of the r -colored suffix of rank i_r

MS with cLCP (2)

Therefore, we only need to sum up cLCP values related to suffixes of s_r versus s_t (as they are computed) for each pair of strings $s_r, s_t \in \mathcal{S}$:



$$\text{Score}(s_r, s_t) = \frac{\sum_{j=1}^{|s_r|} \text{MS}(s_r, s_t)[j]}{|s_r|} = \left(\sum_{\substack{i_r \in [1..M] \\ \text{id}(S)[i_r] = r}} \text{cLCP}[i_r][t] \right) / |s_r|$$

Score Matrix Computation (sketch)

```
Initialize Score[1, m][1, m] = [[0, ..., 0], ..., [0, ..., 0]]
  for t ← 1 to m do
    Dt ← generate_D(t);
    lcpt ← generate_lcp(t);
    score[][t] ← COMPUTECOLUMNSCORE(t, Dt, lcpt, lcp(S), id(S));
COMPUTEACS(score);
```


Score Matrix Computation (sketch)

```
Initialize Score[1, m][1, m] = [[0, ..., 0], ..., [0, ..., 0]]
  for t ← 1 to m do
    Dt ← generate_D(t);
    lcpt ← generate_lcp(t);
    score[ ][t] ← COMPUTECOLUMNSCORE(t, Dt, lcpt, lcp(S), id(S));
COMPUTEACS(score);
```

Each column $score[][t]$ can be computed *separately* from the others, provided D_t and lcp_t are pre-computed.

Score Matrix Computation (sketch)

```
Initialize Score[1, m][1, m] = [[0, ..., 0], ..., [0, ..., 0]]
  for t ← 1 to m do
    Dt ← generate_D(t);
    lcpt ← generate_lcp(t);
    score[ ][t] ← COMPUTECOLUMNSCORE(t, Dt, lcpt, lcp(S), id(S));
COMPUTEACS(score);
```

Straightforward parallel implementation: distinct columns assigned to distinct threads working concurrently.

COMPUTECOLUMNSCORE(t)

```
1: procedure (id( $S$ )[1,  $N$ ], lcp( $S$ )[1,  $N + 1$ ],  $D_t$ [1,  $N$ ],  $lcp_t$ [1,  $|s_t|$ ],  $\alpha[t]$ ,  $\zeta[t]$ )
2:    $\alpha[t] \leftarrow \infty$ 
3:    $\zeta[t] \leftarrow 0$ 
4:    $h_t \leftarrow 1$  ▷ index for scanning  $lcp_t$ 
5:   for  $i \leftarrow 1$  to  $N$  do
6:     if id( $S$ )[ $i$ ]  $\neq t$  then ▷ We are inside a  $t$ -interval:  $[t_1, t_2]$ 
7:        $\alpha[t] \leftarrow \min\{\alpha[t], \text{lcp}[i]\}$ 
8:        $\zeta[t] \leftarrow \max\{\zeta[t], D_t[i] - 1\}$ 
9:       if  $\alpha[t] > lcp_t[i_i]$  then
10:        Score[id[ $i$ ]][ $t$ ]+  $\leftarrow \alpha[t]$ 
11:      else
12:        Score[id[ $i$ ]][ $t$ ]+  $\leftarrow \max\{\alpha[t], \zeta[t], lcp_t[i_i]\}$ 
13:      else ▷ A new  $t$ -interval starts, next  $[t_1, t_2]$ .
14:         $h_t ++$ 
15:         $\alpha[t] \leftarrow \infty$ 
16:         $\zeta[t] \leftarrow 0$ 
```

Preliminary Experiments

- Two collections of genomes (the first one contains 932 genomes and the second one contains 4,938 genomes).
- In both cases, the value $|s_x|$ is greater than the average length of the strings in the respective collection.
- k -Mismatch Average Common Substring approach tool¹ (kmacs) with $k = 0$.

	Size (Gbytes)	Min length	Max length	Max lcp	Program	Wall clock (mm:ss)	Memory (Kbytes)
1	3.434	1,080,084	10,657,107	1,711,194	new cLCP-mACS	2:34	110,412
					cLCP-mACS	13:37	10,716
					kmacs*	23:30	4,213,364
2	9.258	744	14,782,125	5,714,157	new cLCP-mACS	7:43	206,164
					cLCP-mACS	40:21	10,780
					kmacs*	57:43	9,637,964

$|s_x| = 5, 650, 368$ for the first collection and $|s_x| = 3, 571, 103$ for the second one.

All tests were done on a MacBook Pro (13-inch), Intel Core i7 at 3, 5 GHz, with 16 GB of RAM, HDD of type SSD

¹C.-A. Leimeister and M. Burkhard, **Kmacs: the k-mismatch average common substring approach to alignment-free sequence comparison**. *Bioinformatics*, 30(14), 2000-2008.

Future work

- Design a dynamic version of our tool (cLCP can be efficiently auto-updated by removing o inserting strings)
- Solve the many-to-many pairwise ACS problem on a collection of strings or between all strings of a collection versus all strings of another collection
- Use cLCP to define new similarity measures for string collections

Thanks for your attention!