

Edit and Alphabet-Ordering Sensitivity of Lex-Parse

Yuto Nakashima

Kyushu University

Dominik Köppl

University of Yamanashi

Mitsuru Funakoshi

NTT Comm. Sci. Lab.

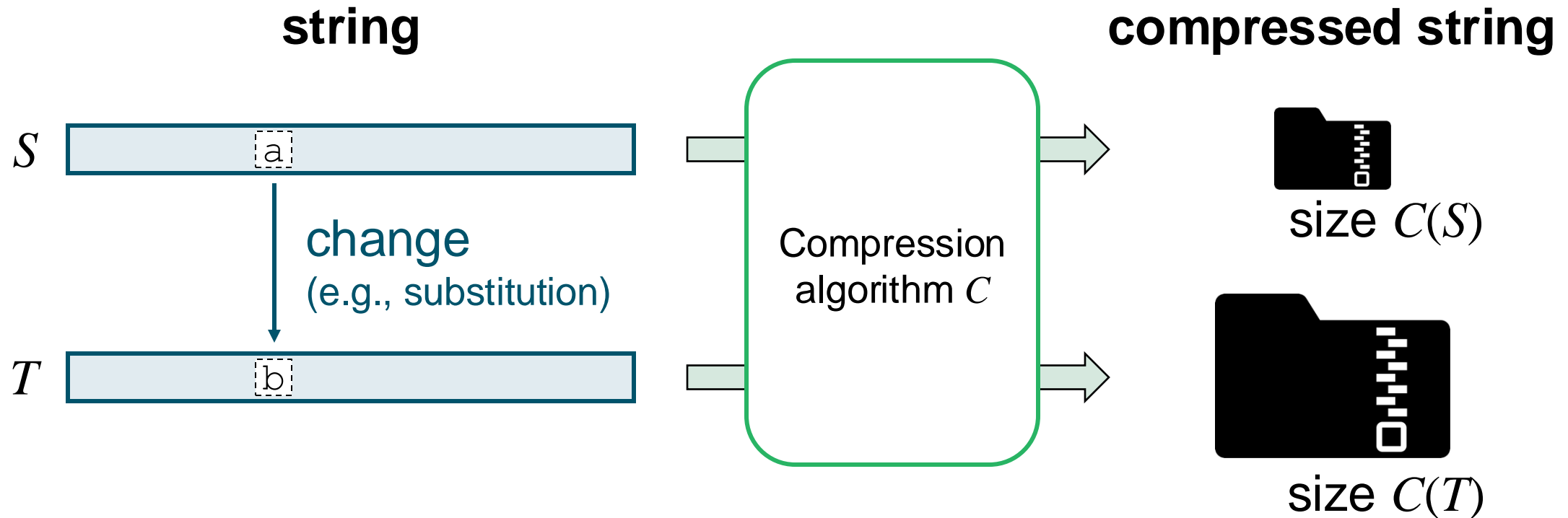
Shunsuke Inenaga

Kyushu University

Hideo Bannai

TMDU

The sensitivity of a string compression algorithm/scheme C is defined as the maximum difference between the size $C(S)$ for a string S and the size $C(T)$ for a single-character edited string T .



How much can the sizes differ in the worst case?

The sensitivity of a string compression algorithm/scheme C is defined as the maximum difference between the size $C(S)$ for a string S and the size $C(T)$ for a single-character edited string T .

There are three edit operations.

- substitution $ab**b**ac \rightarrow a**c**ac$
- insertion $ab-ac \rightarrow ab**c**ac$
- deletion $aba**c** \rightarrow aba-$

In this work,
we study the multiplicative sensitivity.

Multiplicative sensitivity

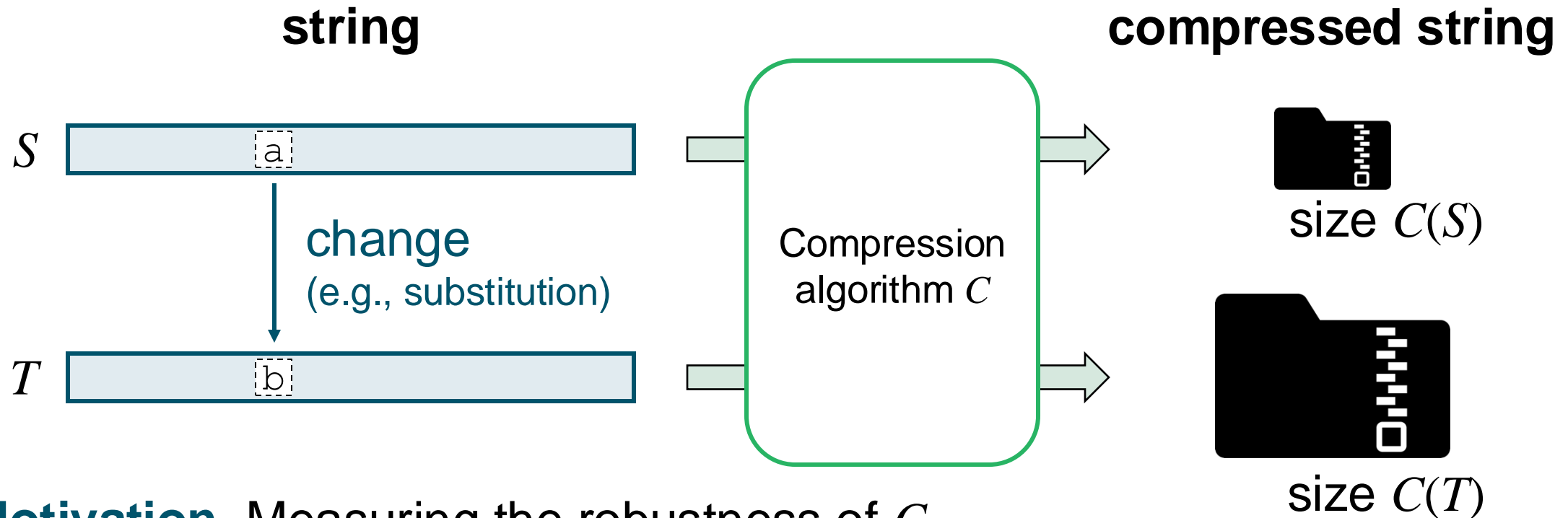
$$\max_{S \in \Sigma^n} \left\{ \frac{C(T)}{C(S)} \mid \text{ED}(S, T) = 1 \right\}$$

edit distance

Additive sensitivity

$$\max_{S \in \Sigma^n} \{ C(T) - C(S) \mid \text{ED}(S, T) = 1 \}$$

The sensitivity of a string compression algorithm/scheme C is defined as the maximum difference between the size $C(S)$ for a string S and the size $C(T)$ for a single-character edited string T .



Motivation. Measuring the robustness of C for errors and/or dynamic changes occurring in the input string.

Multiplicative sensitivity

Measures	Edit operations	Upper bounds	Lower bounds
Smallest bidirectional scheme (b)	all	2	2
Smallest grammar (g^*)	all	2	—
RePair (g_{rpair})	all	$O((n / \log n)^{2/3})$	—
Run-length BWT (r)	all	$O(\log r \log n)$	$\Omega(\log n)$ [Giuliani et al., 2023]
LZ77 (z_{77})	all	2	2
LZSS (z_{ss})	ins.	2	2
	sub./del.	3	3
LZ78 (z_{78})	all	—	$\Omega(n^{1/4})$ [Lagarde & Perifel, 2018]

See the following paper for more results.

“Sensitivity of string compressors and repetitiveness measures”

T. Akagi, M. Funakoshi, S. Inenaga - Information & Computation, 2023

Multiplicative sensitivity

Measures	Edit operations	Upper bounds	Lower bounds
Smallest bidirectional scheme (b)	all	2	2
Smallest grammar (g^*)	all	2	–
RePair (g_{rpair})	all	$O((n / \log n)^{2/3})$	–
Run-length BWT (r)	all	$O(\log r \log n)$	$\Omega(\log n)$ [Giuliani et al., 2023]
LZ77 (z_{77})	all	2	2
LZSS (z_{ss})	ins.	2	2
	sub./del.	3	3
LZ78 (z_{78})	all	–	$\Omega(n^{1/4})$ [Lagarde & Perifel, 2018]
Lex-parse (v)	Open		

Def. The **lex-parse** of a string S is a parsing s_1, \dots, s_v , s.t. each phrase s_j starting at position $i = 1 + \sum_{k < j} |s_k|$ is $S[i..i + \max\{1, \ell\} - 1]$, where ℓ is the length of the longest common prefix between $S[i..]$ and its lex. predecessor $S[i'..]$.

input string

1	2	3	4	5	6	7	8	9
a	a	b	a	b	a	a	b	a

The lex-parse of a string is a string parsing which is defined by the lex. ordering of suffixes of the input string.

We can explain by SA and LCP array.

Suffix array	LCP array	Suffixes sorted in lex. order
9	0	a
6	1	a a b a
1	4	<u>a</u> a b a b a a b a
7	1	a b a
4	3	a b a a b a
2	3	a b a b a a b a
8	0	b a
5	2	b a a b a
3	2	b a b a a b a

Def. The **lex-parse** of a string S is a parsing s_1, \dots, s_v , s.t. each phrase s_j starting at position $i = 1 + \sum_{k < j} |s_k|$ is $S[i..i + \max\{1, \ell\} - 1]$, where ℓ is the length of the longest common prefix between $S[i..]$ and its lex. predecessor $S[i'..]$.

input string

1 2 3 4 5 6 7 8 9
 a a b a b a **a** **b** **a**

Suffix array:

i -th element represents the starting position of the lex. i -th smallest suffix.

LCP array: i -th element represents the length of the longest common prefix between the lex. i -th & $(i-1)$ -th smallest suffix.

Suffix array	LCP array	Suffixes sorted in lex. order
9	0	a
6	1	a a b a
1	4	<u>a</u> a b a b a a b a
7	<u>1</u>	a b a
4	3	a b a a b a
2	3	a b a b a a b a
8	0	b a
5	2	b a a b a
3	2	b a b a a b a

Def. The **lex-parse** of a string S is a parsing s_1, \dots, s_v , s.t. each phrase s_j starting at position $i = 1 + \sum_{k < j} |s_k|$ is $S[i..i + \max\{1, \ell\} - 1]$, where ℓ is the length of the longest common prefix between $S[i..]$ and its lex. predecessor $S[i'..]$.

input string

1 2 3 4 | 5 6 7 8 9
a a b a | b a a b a



Phrase starting at position 1

The length of the LCP between the suffix at position 1 and the its predecessor at position 6 is **4**.

We choose “**aaba**” as the phrase.

Suffix array	LCP array	Suffixes sorted in lex. order
9	0	a
6	1	a a b a
<u>1</u>	4	a a b a b a a b a
7	1	a b a
4	3	a b a a b a
2	3	a b a b a a b a
8	0	b a
5	2	b a a b a
3	2	b a b a a b a

Def. The **lex-parse** of a string S is a parsing s_1, \dots, s_v , s.t. each phrase s_j starting at position $i = 1 + \sum_{k < j} |s_k|$ is $S[i..i + \max\{1, \ell\} - 1]$, where ℓ is the length of the longest common prefix between $S[i..]$ and its lex. predecessor $S[i'..]$.

input string

1 2 3 4 | 5 6 | 7 8 9
 a a b a | **b** **a** | a b a

Phrase starting at position 5

The length of the LCP between the suffix at position 5 and the its pred. at position 8 is **2**.

We choose “**ba**” as the phrase.



Suffix array	LCP array	Suffixes sorted in lex. order
9	0	a
6	1	a a b a
1	4	a a b a b a a b a
7	1	a b a
4	3	a b a a b a
2	3	a b a b a a b a
8	0	b a
<u>5</u>	2	b a a b a
3	2	b a b a a b a

Def. The **lex-parse** of a string S is a parsing s_1, \dots, s_v , s.t. each phrase s_j starting at position $i = 1 + \sum_{k < j} |s_k|$ is $S[i..i + \max\{1, \ell\} - 1]$, where ℓ is the length of the longest common prefix between $S[i..]$ and its lex. predecessor $S[i'..]$.

input string

<u>1</u>	2	3	4	<u>5</u>	6	<u>7</u>	8	9
a	a	b	a	b	a	a	b	a

Phrase starting at position 7

The length of the LCP between the suffix at position 7 and the its pred. at position 1 is **1**.

We choose “**a**” as the phrase.



Suffix array	LCP array	Suffixes sorted in lex. order
9	0	a
6	1	a a b a
1	4	a a b a b a a b a
<u>7</u>	1	a b a
4	3	a b a a b a
2	3	a b a b a a b a
8	0	b a
5	2	b a a b a
3	2	b a b a a b a

Def. The **lex-parse** of a string S is a parsing s_1, \dots, s_v , s.t. each phrase s_j starting at position $i = 1 + \sum_{k < j} |s_k|$ is $S[i..i + \max\{1, \ell\} - 1]$, where ℓ is the length of the longest common prefix between $S[i..]$ and its lex. predecessor $S[i'..]$.

input string

<u>1</u>	2	3	4	<u>5</u>	6	<u>7</u>	<u>8</u>	<u>9</u>
a	a	b	a	b	a	a	b	a

Phrases starting at positions 8 & 9

The length of the LCP is **0**
 (since it is the smallest suffix
 that begins with the character).

We choose the character as the phrase.

Suffix array	LCP array	Suffixes sorted in lex. order
9	0	a
6	1	a a b a
1	4	a a b a b a a b a
7	1	a b a
4	3	a b a a b a
2	3	a b a b a a b a
8	0	b a
5	2	b a a b a
3	2	b a b a a b a

The bidirectional macro scheme is a compression scheme as follows.

Def. [Storer & Szymanski, 1982] A **bidirectional macro scheme** of a string S is a set of directives of the following two types:

- $S[i..j] \leftarrow S[i'..j']$ (i.e. copy $S[i'..j']$ in $S[i..j]$), or
- $S[i] \leftarrow c$, with $c \in \Sigma$ (i.e. assign character c to $S[i]$).

1 2 3 4 5 6 7 8 9
 $S =$ a a b a b a a b a

A BMS of S

$[4, 4] \leftarrow a$

$[5, 5] \leftarrow b$

$[1, 3] \leftarrow [6, 8]$

$[7, 9] \leftarrow [4, 6]$

$[6, 6] \leftarrow [4, 4]$

The bidirectional macro scheme is a compression scheme as follows.

Def. [Storer & Szymanski, 1982] A **bidirectional macro scheme** of a string S is a set of directives of the following two types:

- $S[i..j] \leftarrow S[i'..j']$ (i.e. copy $S[i'..j']$ in $S[i..j]$), or
- $S[i] \leftarrow c$, with $c \in \Sigma$ (i.e. assign character c to $S[i]$).

1 2 3 4 5 6 7 8 9
 $S =$ a a b a b a a b a

Decoding

1 2 3 4 5 6 7 8 9
a b

A BMS of S

$[4, 4] \leftarrow a$

$[5, 5] \leftarrow b$

$[1, 3] \leftarrow [6, 8]$

$[7, 9] \leftarrow [4, 6]$

$[6, 6] \leftarrow [4, 4]$

The bidirectional macro scheme is a compression scheme as follows.

Def. [Storer & Szymanski, 1982] A **bidirectional macro scheme** of a string S is a set of directives of the following two types:

- $S[i..j] \leftarrow S[i'..j']$ (i.e. copy $S[i'..j']$ in $S[i..j]$), or
- $S[i] \leftarrow c$, with $c \in \Sigma$ (i.e. assign character c to $S[i]$).

1 2 3 4 5 6 7 8 9
 $S =$ a a b a b a a b a

Decoding

1 2 3 4 5 6 7 8 9
a b a

A BMS of S

$[4, 4] \leftarrow a$ $[5, 5] \leftarrow b$
 $[1, 3] \leftarrow [6, 8]$ $[7, 9] \leftarrow [4, 6]$
 $[6, 6] \leftarrow [4, 4]$

The bidirectional macro scheme is a compression scheme as follows.

Def. [Storer & Szymanski, 1982] A **bidirectional macro scheme** of a string S is a set of directives of the following two types:

- $S[i..j] \leftarrow S[i'..j']$ (i.e. copy $S[i'..j']$ in $S[i..j]$), or
- $S[i] \leftarrow c$, with $c \in \Sigma$ (i.e. assign character c to $S[i]$).

1 2 3 4 5 6 7 8 9
 $S =$ a a b a b a a b a

Decoding

1 2 3 4 5 6 7 8 9
a b a a b a

A BMS of S

$[4, 4] \leftarrow a$ $[5, 5] \leftarrow b$
 $[1, 3] \leftarrow [6, 8]$ $[7, 9] \leftarrow [4, 6]$
 $[6, 6] \leftarrow [4, 4]$

The bidirectional macro scheme is a compression scheme as follows.

Def. [Storer & Szymanski, 1982] A **bidirectional macro scheme** of a string S is a set of directives of the following two types:

- $S[i..j] \leftarrow S[i'..j']$ (i.e. copy $S[i'..j']$ in $S[i..j]$), or
- $S[i] \leftarrow c$, with $c \in \Sigma$ (i.e. assign character c to $S[i]$).

	1	2	3	4	5	6	7	8	9
$S =$	a	a	b	a	b	a	a	b	a

Decoding

1	2	3	4	5	6	7	8	9
a	a	b	a	b	<u>a</u>	<u>a</u>	<u>b</u>	a

A BMS of S

$[4, 4] \leftarrow a$	$[5, 5] \leftarrow b$
$[1, 3] \leftarrow [6, 8]$	$[7, 9] \leftarrow [4, 6]$
$[6, 6] \leftarrow [4, 4]$	

- The lex-parse gives a bidirectional macro scheme.
 - ▷ LCP means that there is another occ. in the text.
 - ▷ The position that has no common prefix is a **single character phrase**.

1 2 3 4 | 5 6 | 7 | 8 | 9
a a b a | b a | a | b | a



Suffix array	LCP array	Suffixes sorted in lex. order
9	0	a
6	1	a a b a
1	4	a a b a b a a b a
7	1	a b a
4	3	a b a a b a
2	3	a b a b a a b a
8	0	b a
5	2	b a a b a
3	2	b a b a a b a

BMS by lex-parse

- [9, 9] ← a
- [8, 8] ← b
- [1, 4] ← [6, 9]
- [5, 6] ← [8, 9]
- [7, 7] ← [1, 1]

The lex-parse gives a bidirectional macro scheme.

Def. [Storer & Szymanski, 1982] A **bidirectional macro scheme** of a string S is a set of directives of the following two types:

- $S[i..j] \leftarrow S[i'..j']$ (i.e. copy $S[i'..j']$ in $S[i..j]$), or
- $S[i] \leftarrow c$, with $c \in \Sigma$ (i.e. assign character c to $S[i]$).

- $b(S)$ denotes the minimum number of directives in BMS of a string S .
- $v(S)$ denotes the size of the lex-parse.
- For any string S , $b(S) \leq v(S)$ holds [Navarro et al., 2021].

Lex-parse gives a BMS.

- Lex-parse was proposed as **a new variant in a family of ordered parsings** that is considered as a generalization of the LZ parsing and a subset of bidirectional macro schemes.
 - ▷ LZ parsing: Each phrase refers **a smaller suffix w.r.t. the text position.**
 - ▷ Lex-parse: Each phrase refers **a smaller suffix w.r.t. the lex. order.**
- Combinatorial studies on lex-parse can lead us to further understanding in string repetitiveness measures and compressors.
 - ▷ A direct relation $\nu \in O(r)$ between the size of the lex-parse ν and the size r of the RLBWT, one of the most important dictionary compressors, holds [Navarro et al., 2021].
 - ▷ A direct relation between the lex-parse ν and the LZ parsing z is open.

Multiplicative edit-sensitivity

Measures	Edit operations	Upper bounds	Lower bounds
Lex-parse (v)	all	$O(\log(n/b))$	$\Omega(\log n)$

Since we use a family of strings s.t. $b = O(1)$, this bound does not contradict the upper bound.

Overview

- The upper bound can be obtained by combining some known results (regarding the lex-parse and the bidirectional macro scheme).
- The lower bound can be obtained by the **Fibonacci word**.
 - ▷ Combinatorial properties on the Fibonacci word
 - ▷ Characterizing the lex-parse by **Lyndon factorization**

We introduce a new sensitivity of compressors by alphabet orderings.

- If a compressor C is defined by using the lexicographic order, the size $C(S)$ depends on alphabet orderings.

Multiplicative sensitivity

$$\max_{S \in \Sigma^n} \left\{ \frac{C(S, \prec_2)}{C(S, \prec_1)} \mid \prec_1, \prec_2 \in A \right\}$$

Additive sensitivity

$$\max_{S \in \Sigma^n} \{ C(S, \prec_2) - C(S, \prec_1) \mid \prec_1, \prec_2 \in A \}$$

For instance,

- **Lex-parse**
- (Run-length) Burrows-Wheeler transform
- GCIS (a grammar comp. algorithm based on induced sorting)

In this work, we investigate the multiplicative sensitivity.

Optimization problems of these kinds of structures have been studied.

■ Complexity (NP-hardness)

- ▷ Minimization for the RLBWT [Bentley et al., 2020]
- ▷ Minimization/Maximization for the Lyndon factorization [Gibney & Thankachan, 2021]
- ▷ Minimization for the minimizer [Verbeek et al., 2024]

■ Exact algorithms/heuristics

- ▷ for RLBWT and Lyndon factorization [Cenzato et al., 2023] [Clare & Daykin, 2019] [Clare et al., 2019] [Major et al., 2020]

Multiplicative AO-sensitivity

Measures	Upper bounds	Lower bounds	
Run-length BWT (r)	$O(\log^2 n)$	$\Omega(\log n)$	from known results
Lyndon factorization [†]	$O(n)$	$\Omega(n)$	obvious
Lex-parse	$O(\log(n/b))$	$\Omega(\log n)$	our paper

Overview

$$b = O(1)$$

[†] Lyndon factorization is not a compressor, but a popular object depending on the lex. order.

- The upper bound can be obtained by combining some known results (regarding the lex-parse and the bidirectional macro scheme).
- The lower bound can be obtained by the Fibonacci word over $\{a, b\}$.
 - ▷ Although the results for “ $a < b$ ” have been proven [Navarro et al., 2021], we give alternative proofs for this case that leads us to the proof for the case when “ $b < a$ ”.

Multiplicative **Edit**-sensitivity (for all operations)

Measures	Upper bounds	Lower bounds
Lex-parse (v)	$O(\log (n/b))$	$\Omega(\log n)$

Multiplicative **AO**-sensitivity

Measures	Upper bounds	Lower bounds
Lex-parse (v)	$O(\log (n/b))$	$\Omega(\log n)$

In this talk, I will explain

- the upper bounds for both sensitivities, and
- the lower bound for the edit-sensitivity.

The upper bounds can be obtained by combining some known results.

Lemma. For any string S and T such that $ED(S, T) = 1$,

- $b(S) \leq v(S)$ [Navarro et al., 2021],
- $v(S) \in O(b(S) \log(n/b(S)))$ [Navarro et al., 2021],
- $b(T) \leq 2b(S)$ [Akagi et al., 2023].

■ Multiplicative Edit-sensitivity

$$\frac{v(T)}{v(S)} \leq \frac{v(T)}{b(S)} \in O\left(\frac{b(T) \log \frac{n}{b(T)}}{b(S)}\right) \subseteq O\left(\frac{b(S) \log \frac{n}{b(S)}}{b(S)}\right) = O\left(\log \frac{n}{b(S)}\right)$$

■ Multiplicative AO-sensitivity

$$\frac{v(S, <_2)}{v(S, <_1)} \leq \frac{v(S, <_2)}{b(S)} \in O\left(\frac{b(S) \log \frac{n}{b(S)}}{b(S)}\right) = O\left(\log \frac{n}{b(S)}\right)$$

In the rest of this talk,
I will explain our ideas for the lower bound for edit-sensitivity.

The formal statement for our result is given as follows.

Theorem. There exists a family of strings α_k, β_k that satisfies $\text{ED}(\alpha_k, \beta_k) = 1$ and $v(\beta_k) / v(\alpha_k) \in \Omega(\log n)$.

- We use (finite) **Fibonacci words** F_{2k} as the base strings α_k .

Def. The i -th Fibonacci word over $\{a, b\}$ is defined as follows:

$$F_1 = b, F_2 = a, F_i = F_{i-1} \cdot F_{i-2}.$$

f_i	F_i
1	$F_1 = b$
1	$F_2 = a$
2	$F_3 = ab$
3	$F_4 = aba$
5	$F_5 = abaab$
8	$F_6 = \text{abaababa}$
13	$F_7 = \text{abaababaabaab}$
21	$F_8 = \text{abaababaabaababaababa}$
34	$F_9 = \text{abaababaabaababaababaabaababaabaab}$
55	$F_{10} = \text{abaababaabaababaababaabaababaabaababaabaababaabaababaabaababaabaab}$

f_i denotes the length $|F_i|$.
 (= f_i corresponds the Fibonacci sequence.)

Def. The i -th Fibonacci word over $\{a, b\}$ is defined as follows:

$$F_1 = b, F_2 = a, F_i = F_{i-1} \cdot F_{i-2}.$$

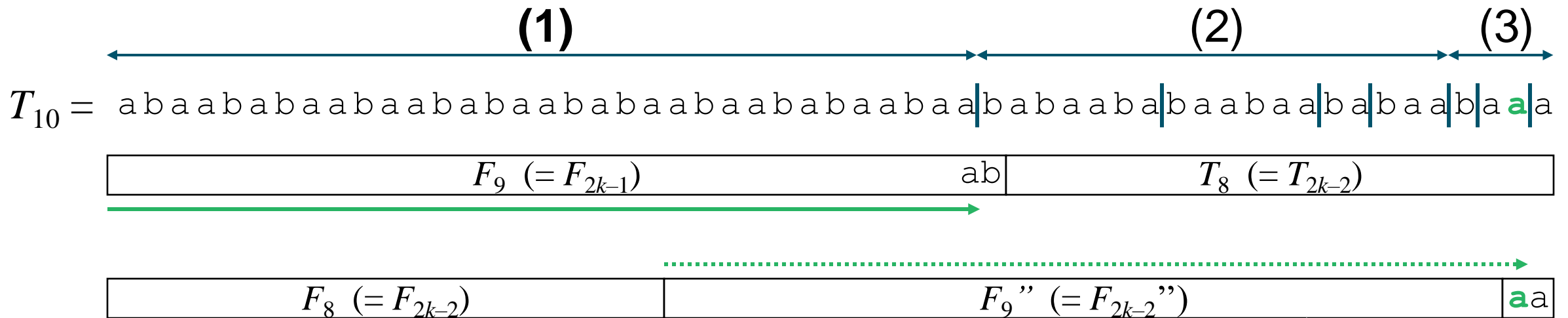
f_i	F_i	
1	F_1	= b
1	F_2	= a
2	F_3	= ab
3	F_4	= a ba
5	F_5	= aba ab
8	F_6	= abaaba ba
13	F_7	= abaababaaba ab
21	F_8	= abaababaabaababaaba ba
34	F_9	= abaababaabaababaabaababaabaaba ab
55	F_{10}	= abaababaabaababaabaababaabaabaababaabaabaabaabaabaabaabaabaabaaba ba

Property. For every **even** (resp. odd) i , the suffix of length 2 is **ba** (resp. ab).

Property. **aaa** and **bb** cannot appear.

There are three types of phrases:

(1) first phrase, (2) inductive phrases, and (3) last three phrases.



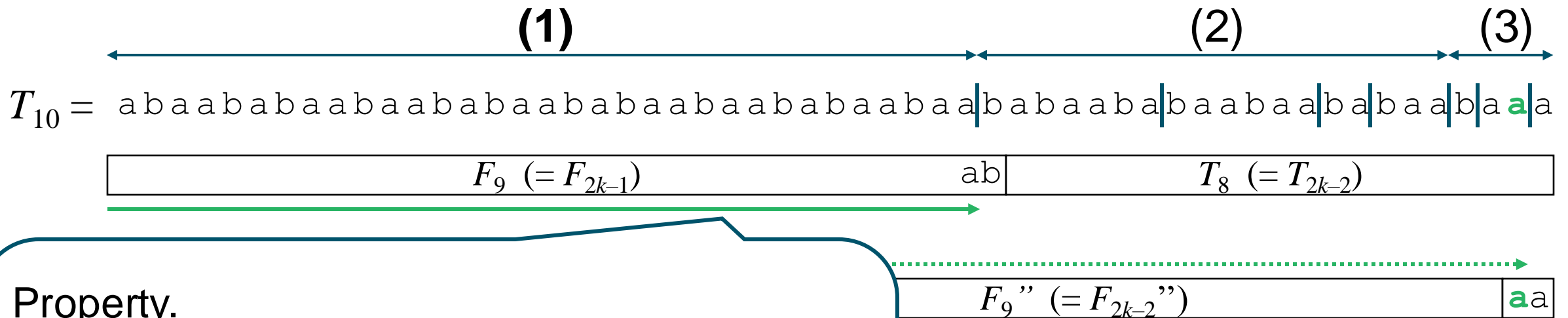
■ We can show the reference by using the following properties.

- ▷ F_{i-1} only occurs as a prefix of F_i .
- ▷ Suffix "aaa" cannot occur as an infix of T_i .

For any string S ,
 $S' = S[1..|S|-1]$, $S'' = S[1..|S|-2]$.

There are three types of phrases:

(1) first phrase, (2) inductive phrases, and (3) last three phrases.



Property.

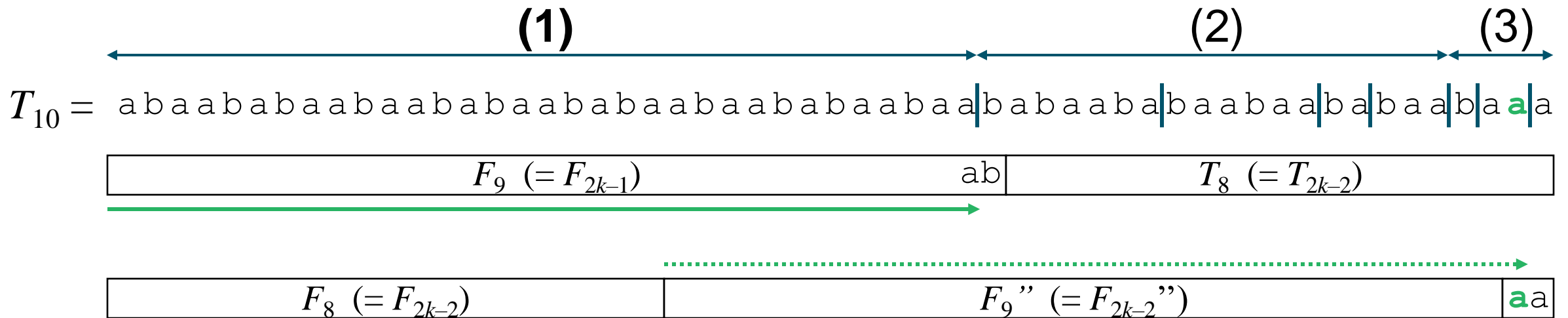


flip $ab \leftrightarrow ba$ T_i .

following properties.

There are three types of phrases:

(1) first phrase, (2) inductive phrases, and (3) last three phrases.



■ We can show the reference by using the following properties.

- ▷ F_{i-1} only occurs as a prefix of F_i .
- ▷ Suffix "aaa" cannot occur as an infix of T_i .

Def. A string λ is a Lyndon word, if λ is lexicographically smaller than any of its non-empty proper suffixes.

		proper suffixes
$\lambda = a b a b b$	\prec	$b a b b$
	\prec	$a b b$
	\prec	$b b$
	\prec	b

$ababb$ is a Lyndon word.

Def. A factorization $\lambda_1^{p_1}, \dots, \lambda_m^{p_m}$ of S is the Lyndon factorization $LF(S)$, if

- $\lambda_1 > \dots > \lambda_m$ are Lyndon words, and
- $p_i \geq 1$ for all $1 \leq i \leq m$.

$$\lambda = abb|abb|ababb|ab|ab|ab|a$$

$$abb = abb > ababb > ab = ab = ab > a = a$$

$$LF(\lambda) = (abb)^2 \quad ababb \quad (ab)^3 \quad (a)^2$$

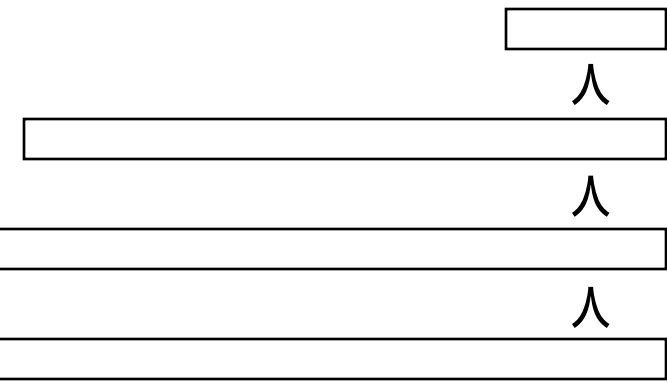
For any string S , the Lyndon factorization of S is unique.

Def. A factorization $\lambda_1^{p_1}, \dots, \lambda_m^{p_m}$ of S is the Lyndon factorization $LF(S)$, if

- $\lambda_1 \succ \dots \succ \lambda_m$ are Lyndon words, and
- $p_i \geq 1$ for all $1 \leq i \leq m$.

$$\lambda = a b b \dot{a} b b | a b a b b | a b \dot{a} b \dot{a} b | a \dot{a} a$$

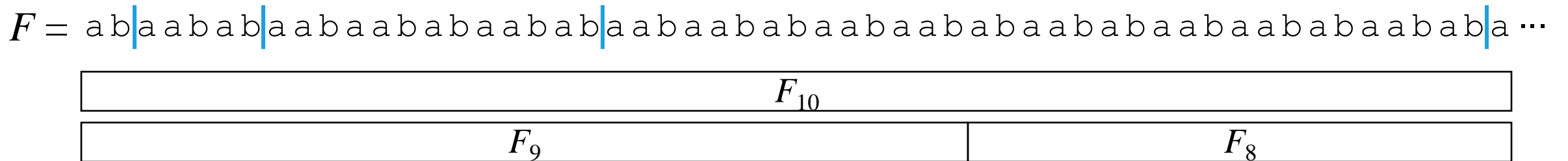
Suffixes that are concatenation of Lyndon factors are sorted in lex. order.



Property. For any integer $i \in [1, m-1]$, $\lambda_i^{p_i} \dots \lambda_m^{p_m} \succ \lambda_{i+1}^{p_{i+1}} \dots \lambda_m^{p_m}$ holds.

Def. Let $F = \lim_{i \rightarrow \infty} F_i$ be the infinite Fibonacci word, and $LF(F) = \ell_1, \ell_2, \dots$ be the (infinite) Lyndon factorization of F .

Lemma. [Melançon, 2000] Let ϕ be a string morphism s.t. $\phi(a) = aab, \phi(b) = ab$. $\ell_1 = ab, \ell_{i+1} = \phi(\ell_i)$, and $|\ell_i| = f_{2i+1}$ holds.



To show our result, we consider the Lyndon factorization of $F_i'' = T_i''$ (the string that can be obtained by removing the edited suffix “aa”).

Multiplicative **Edit**-sensitivity (for all operations)

Measures	Upper bounds	Lower bounds
Lex-parse (v)	$O(\log (n/b))$	$\Omega(\log n)$

Multiplicative **AO**-sensitivity

Measures	Upper bounds	Lower bounds
Lex-parse (v)	$O(\log (n/b))$	$\Omega(\log n)$

Further work (on alphabet-orderings)

- Problem of computing optimal alphabet orderings for the lex-parse.
 - The problems for the RLBWT [Bentley et al., 2020] and the Lyndon factorization [Gibney & Thankachan, 2021] are known to be NP-hard.