# On the approximation ratio of LZ-End to LZ77

Takumi Ideue, Takuya Mieno, Mitsuru Funakoshi,
Yuto Nakashima, Shunsuke Inenaga, Masayuki Takeda
Kyushu University, Japan

# LZ77 vs LZ-End

LZ77 [Ziv and Lempel, 1977] is the smallest greedy parsing allowing for left-to-right (de)compression.

LZ-End [Kreft and Navarro, 2013] is an LZ77-like parsing allowing for fast substring extraction,
but <u>the number of its phrases is larger than that of LZ77</u>.

**Theorem:** [This work]

There exist **binary** strings $S$ such that:

$$\frac{z_{End}(S)}{z_{77}(S)} \to 2 \ (|S| \to \infty).$$

$z_{End}(S)$: # of LZ-End phrases of $S$
$z_{77}(S)$: # of LZ77 phrases of $S$

**Definition:**

The non-overlapping Lempel-Ziv 77 factorization (LZ77) of a string $T$ is the factorization $LZ_{77}(T) = p_1, \ldots, p_z$ of $T$ such that: Each phrase $p_i$ ($1 \leq i \leq z - 1$) satisfies the following condition.

- $p_i[1, |p_i| - 1]$ is <u>the longest prefix</u> of $p_i \cdots p_z$ which occurs in $p_1 \cdots p_{i-1}$.

The last phrase $p_z$ can be **a suffix** of $T$ which occurs in $p_1 \cdots p_{i-1}$.

$z$ is the number of phrases

E.g.)

$$
\begin{array}{c}
\quad 1\ 2\ \ 3\ \ 4\ \ 5\ \ 6\ \ 7\ \ 8\ \ 9\ 10\ 11\ 12\ 13\ 14\ 15\ 16\ 17 \\
LZ_{77}(T) = \text{a b a a b a b b a b b a b a a b b}
\end{array}
$$

# LZ77 [Ziv and Lempel, 1977]

**Definition:**

The non-overlapping Lempel-Ziv 77 factorization (LZ77) of a string $T$ is the factorization $\mathrm{LZ}_{77}(T) = p_1, \ldots, p_z$ of $T$ such that: Each phrase $p_i$ $(1 \le i \le z - 1)$ satisfies the following condition.

- $p_i[1, |p_i| - 1]$ is <u>the longest prefix</u> of $p_i \cdots p_z$ which occurs in $p_1 \cdots p_{i-1}$.

The last phrase $p_z$ can be **a suffix** of $T$ which occurs in $p_1 \cdots p_{i-1}$.

$z$ is the number of phrases

E.g.)

$$\mathrm{LZ}_{77}(T) = \begin{array}{cccccccccccccccccc} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 & 13 & 14 & 15 & 16 & 17 \\ \mathrm{a} & \mathrm{b} & \mathrm{a} & \mathrm{a} & \mathrm{b} & \mathrm{a} & \mathrm{b} & \mathrm{b} & \mathrm{a} & \mathrm{b} & \mathrm{b} & \mathrm{a} & \mathrm{b} & \mathrm{a} & \mathrm{a} & \mathrm{b} & \mathrm{b} \end{array}$$

First occurrence

**Definition:**

The non-overlapping Lempel-Ziv 77 factorization (LZ77) of a string $T$ is the factorization $LZ_{77}(T) = p_1, \ldots, p_z$ of $T$ such that: Each phrase $p_i$ ($1 \leq i \leq z - 1$) satisfies the following condition.

- $p_i[1, |p_i| - 1]$ is <u>the longest prefix</u> of $p_i \cdots p_z$ which occurs in $p_1 \cdots p_{i-1}$.

The last phrase $p_z$ can be **a suffix** of $T$ which occurs in $p_1 \cdots p_{i-1}$.

$z$ is the number of phrases

E.g.)

$$\begin{array}{cccccccccccccccccc} & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 & 13 & 14 & 15 & 16 & 17 \end{array}$$

$$LZ_{77}(T) = \text{a} | \text{b} | \text{a a b a b b a b b a b a a b b}$$

First occurrence

# LZ77 [Ziv and Lempel, 1977]

**Definition:**

The non-overlapping Lempel-Ziv 77 factorization (LZ77) of a string $T$ is the factorization $\mathrm{LZ}_{77}(T) = p_1, \ldots, p_z$ of $T$ such that: Each phrase $p_i$ ($1 \leq i \leq z - 1$) satisfies the following condition.

- $p_i[1, |p_i| - 1]$ is <u>the longest prefix</u> of $p_i \cdots p_z$ which occurs in $p_1 \cdots p_{i-1}$.

The last phrase $p_z$ can be **a suffix** of $T$ which occurs in $p_1 \cdots p_{i-1}$.

$z$ is the number of phrases

E.g.)

$$\mathrm{LZ}_{77}(T) = \underset{1}{\underline{a}}|\underset{2}{b}|\underset{3}{\underline{a}}\,\underset{4}{a}|\underset{5}{b}\,\underset{6}{a}\,\underset{7}{b}\,\underset{8}{b}\,\underset{9}{a}\,\underset{10}{b}\,\underset{11}{b}\,\underset{12}{a}\,\underset{13}{b}\,\underset{14}{a}\,\underset{15}{a}\,\underset{16}{b}\,\underset{17}{b}$$

The longest prefix of $p_3 \cdots$

6

**Definition:**

The non-overlapping Lempel-Ziv 77 factorization (LZ77) of a string $T$ is the factorization $LZ_{77}(T) = p_1, \ldots, p_z$ of $T$ such that: Each phrase $p_i$ $(1 \leq i \leq z - 1)$ satisfies the following condition.

- $p_i[1, |p_i| - 1]$ is <u>the longest prefix</u> of $p_i \cdots p_z$ which occurs in $p_1 \cdots p_{i-1}$.

The last phrase $p_z$ can be **a suffix** of $T$ which occurs in $p_1 \cdots p_{i-1}$.

$z$ is the number of phrases

E.g.)

$$\begin{array}{cccccccccccccccccc} & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 & 13 & 14 & 15 & 16 & 17 \end{array}$$

$$LZ_{77}(T) = \text{a}|\text{b}|\text{a a}|\text{b a b}|\text{b a b b a b a a b b}$$

The longest prefix of $p_4 \cdots$

# LZ77 [Ziv and Lempel, 1977]

**Definition:**

The non-overlapping Lempel-Ziv 77 factorization (LZ77) of a string $T$ is the factorization $\mathrm{LZ}_{77}(T) = p_1, \ldots, p_z$ of $T$ such that: Each phrase $p_i$ ($1 \leq i \leq z-1$) satisfies the following condition.

- $p_i[1, |p_i|-1]$ is <u>the longest prefix</u> of $p_i \cdots p_z$ which occurs in $p_1 \cdots p_{i-1}$.

The last phrase $p_z$ can be **a suffix** of $T$ which occurs in $p_1 \cdots p_{i-1}$.

> $z$ is the number of phrases

E.g.)

$$\mathrm{LZ}_{77}(T) = \underset{\substack{1\ \ 2\ \ 3\ \ 4\ \ 5\ \ 6\ \ 7\ \ 8\ \ 9\ \ 10\ 11\ 12\ 13\ 14\ 15\ 16\ 17}}{\mathrm{a}|\mathrm{b}|\mathrm{a}\ \mathrm{a}|\mathrm{b}\ \mathrm{a}\ \mathrm{b}|\mathrm{b}\ \mathrm{a}\ \mathrm{b}\ \mathrm{a}\ \mathrm{b}\ \mathrm{a}|\mathrm{a}\ \mathrm{b}\ \mathrm{b}}$$

**Definition:**

The non-overlapping Lempel-Ziv 77 factorization (LZ77) of a string $T$ is the factorization $LZ_{77}(T) = p_1, \ldots, p_z$ of $T$ such that: Each phrase $p_i$ ($1 \leq i \leq z - 1$) satisfies the following condition.

- $p_i[1, |p_i| - 1]$ is <u>the longest prefix</u> of $p_i \cdots p_z$ which occurs in $p_1 \cdots p_{i-1}$.

The last phrase $p_z$ can be **a suffix** of $T$ which occurs in $p_1 \cdots p_{i-1}$.

$z$ is the number of phrases

E.g.)

$$\begin{array}{cccccccccccccccccc} & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 & 13 & 14 & 15 & 16 & 17 \end{array}$$

$$LZ_{77}(T) = a|b|a\ a|b\ a\ b|b\ a\ b\ b\ a\ b\ a|a\ b\ b$$

**Non-overlapping**

**Definition:**

The non-overlapping Lempel-Ziv 77 factorization (LZ77) of a string $T$ is the factorization $LZ_{77}(T) = p_1, …, p_z$ of $T$ such that: Each phrase $p_i$ $(1 \le i \le z - 1)$ satisfies the following condition.

- $p_i[1, |p_i| - 1]$ is <u>the longest prefix</u> of $p_i \cdots p_z$ which occurs in $p_1 \cdots p_{i-1}$.

The last phrase $p_z$ can be **a suffix** of $T$ which occurs in $p_1 \cdots p_{i-1}$.

$z$ is the number of phrases

E.g.)

$$1 \quad 2 \quad 3 \quad 4 \quad 5 \quad 6 \quad 7 \quad 8 \quad 9 \quad 10 \; 11 \; 12 \; 13 \; 14 \; 15 \; 16 \; 17$$

$$LZ_{77}(T) = \text{a} | \text{b} | \text{a a} | \text{b a b} | \text{b a b b} | \text{a b a a b b}$$

The longest prefix of $p_5 \cdots$

**Definition:**

The non-overlapping Lempel-Ziv 77 factorization (LZ77) of a string $T$ is the factorization $LZ_{77}(T) = p_1, \ldots, p_z$ of $T$ such that: Each phrase $p_i$ ($1 \le i \le z - 1$) satisfies the following condition.

- $p_i[1, |p_i| - 1]$ is <u>the longest prefix</u> of $p_i \cdots p_z$ which occurs in $p_1 \cdots p_{i-1}$.

The last phrase $p_z$ can be **a suffix** of $T$ which occurs in $p_1 \cdots p_{i-1}$.

$z$ is the number of phrases

E.g.)

$$
\begin{array}{cccccccccccccccccc}
 & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 & 13 & 14 & 15 & 16 & 17 \\
\end{array}
$$

$$LZ_{77}(T) = \text{a|b|a a|b a b|b a b b|a b a a b b|}$$

The longest prefix of $p_6$

**Definition:**

The non-overlapping Lempel-Ziv 77 factorization (LZ77) of a string $T$ is the factorization $LZ_{77}(T) = p_1, \ldots, p_z$ of $T$ such that: Each phrase $p_i$ $(1 \le i \le z - 1)$ satisfies the following condition.

- $p_i[1, |p_i| - 1]$ is <u>the longest prefix</u> of $p_i \cdots p_z$ which occurs in $p_1 \cdots p_{i-1}$.

The last phrase $p_z$ can be **a suffix** of $T$ which occurs in $p_1 \cdots p_{i-1}$.

> $z$ is the number of phrases

E.g.)

$$
\begin{array}{cccccccccccccccccc}
 & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 & 13 & 14 & 15 & 16 & 17
\end{array}
$$

$$LZ_{77}(T) = a|b|a\ a|b\ a\ b|b\ a\ b\ b|a\ b\ a\ a\ b\ b|$$

$$p_1\ p_2\ \ p_3\ \ \ \ \ p_4\ \ \ \ \ \ \ \ p_5\ \ \ \ \ \ \ \ \ \ \ \ p_6$$

$$z_{77}(T) = 6$$

> The number of the LZ77 phrases of $T$

12

# LZ-End [Kreft and Navarro, 2013]

**Definition:**

The LZ-End factorization of a string $T$ is the factorization $\mathrm{LZ}_{\mathrm{End}}(T) = q_1, \ldots, q_{z'}$ of $T$ such that:

Each phrase $q_i$ $(1 \leq i \leq z' - 1)$ satisfies the following condition.

- $q_i[1, |q_i| - 1]$ is <u>the longest prefix</u> of $q_i \cdots q_{z'}$ which occurs as a suffix of $q_1 \cdots q_j$ for some $j < i$.

The last phrase $q_{z'}$ can be **a suffix** of $T$ which occurs as a suffix of $q_1 \cdots q_j$ for some $j < z'$.

$z'$ is the number of phrases

E.g.)

$$
\begin{array}{ccccccccccccccccc}
1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 & 13 & 14 & 15 & 16 & 17
\end{array}
$$

$$
\mathrm{LZ}_{\mathrm{End}}(T) = \text{a b a a b a b b a b b a b a a b b}
$$

13

**Definition:**

The LZ-End factorization of a string $T$ is the factorization $\text{LZ}_{\text{End}}(T) = q_1, \ldots, q_{z'}$ of $T$ such that:

$z'$ is the number of phrases

Each phrase $q_i$ ($1 \leq i \leq z' - 1$) satisfies the following condition.

- $q_i[1, |q_i| - 1]$ is <u>the longest prefix</u> of $q_i \cdots q_{z'}$ which occurs as a suffix of $q_1 \cdots q_j$ for some $j < i$.

The last phrase $q_{z'}$ can be **a suffix** of $T$ which occurs as a suffix of $q_1 \cdots q_j$ for some $j < z'$.

E.g.)

$$\text{LZ}_{\text{End}}(T) = \begin{array}{ccccccccccccccccc} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 & 13 & 14 & 15 & 16 & 17 \\ a & b & a & a & b & a & b & b & a & b & b & a & b & a & a & b & b \end{array}$$

First occurrence

14

**Definition:**

The LZ-End factorization of a string $T$ is the factorization $LZ_{End}(T) = q_1, \ldots, q_{z'}$ of $T$ such that:

> $z'$ is the number of phrases

Each phrase $q_i$ $(1 \leq i \leq z' - 1)$ satisfies the following condition.

- $q_i[1, |q_i| - 1]$ is <u>the longest prefix</u> of $q_i \cdots q_{z'}$ which occurs as a suffix of $q_1 \cdots q_j$ for some $j < i$.

The last phrase $q_{z'}$ can be **a suffix** of $T$ which occurs as a suffix of $q_1 \cdots q_j$ for some $j < z'$.

E.g.)

$$\begin{array}{cccccccccccccccccc} & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 & 13 & 14 & 15 & 16 & 17 \end{array}$$

$$LZ_{End}(T) = \text{a} | \text{b} | \text{a a b a b b a b b a b a a b b}$$

> First occurrence

**Definition:**

The LZ-End factorization of a string $T$ is the factorization $\mathrm{LZ}_{\mathrm{End}}(T) = q_1, \ldots, q_{z'}$ of $T$ such that:

> $z'$ is the number of phrases

Each phrase $q_i$ $(1 \le i \le z' - 1)$ satisfies the following condition.

- $q_i[1, |q_i| - 1]$ is <u>the longest prefix</u> of $q_i \cdots q_{z'}$ which occurs as a suffix of $q_1 \cdots q_j$ for some $j < i$.

The last phrase $q_{z'}$ can be **a suffix** of $T$ which occurs as a suffix of $q_1 \cdots q_j$ for some $j < z'$.

E.g.)

```
              1  2  3  4  5  6  7  8  9 10 11 12 13 14 15 16 17
LZ_End(T) =   a| b| a  a| b  a  b  b  a  b  b  a  b  a  a  b  b
```

Suffix of $q_1$

The longest prefix of $q_3$ $\cdots$

16

**Definition:**

The LZ-End factorization of a string $T$ is the factorization $\mathrm{LZ}_{\mathrm{End}}(T) = q_1, \ldots, q_{z'}$ of $T$ such that:

Each phrase $q_i$ $(1 \le i \le z' - 1)$ satisfies the following condition.

- $q_i[1, |q_i| - 1]$ is <u>the longest prefix</u> of $q_i \cdots q_{z'}$ which occurs as a suffix of $q_1 \cdots q_j$ for some $j < i$.

The last phrase $q_{z'}$ can be **a suffix** of $T$ which occurs as a suffix of $q_1 \cdots q_j$ for some $j < z'$.

> $z'$ is the number of phrases

E.g.)

$$
\begin{array}{c}
\phantom{\mathrm{LZ}_{\mathrm{End}}(T) = } \; 1 \;\; 2 \;\; 3 \;\; 4 \;\; 5 \;\; 6 \;\; 7 \;\; 8 \;\; 9 \;\; 10 \;\; 11 \;\; 12 \;\; 13 \;\; 14 \;\; 15 \;\; 16 \;\; 17 \\
\mathrm{LZ}_{\mathrm{End}}(T) = \mathrm{a}\,|\,\mathrm{b}\,|\,\mathrm{a}\;\mathrm{a}\,|\,\mathrm{b}\;\mathrm{a}\,|\,\mathrm{b}\;\mathrm{b}\;\mathrm{a}\;\mathrm{b}\;\mathrm{b}\;\mathrm{a}\;\mathrm{b}\;\mathrm{a}\;\mathrm{a}\;\mathrm{b}\;\mathrm{b}
\end{array}
$$

Suffix of $q_1 q_2$

The longest prefix of $q_4 \cdots$

**Definition:**

The LZ-End factorization of a string $T$ is the factorization $\mathrm{LZ}_{\mathrm{End}}(T) = q_1, \ldots, q_{z'}$ of $T$ such that:

$z'$ is the number of phrases

Each phrase $q_i$ $(1 \leq i \leq z' - 1)$ satisfies the following condition.

- $q_i[1, |q_i| - 1]$ is <u>the longest prefix</u> of $q_i \cdots q_{z'}$ which occurs as a suffix of $q_1 \cdots q_j$ for some $j < i$.

The last phrase $q_{z'}$ can be **a suffix** of $T$ which occurs as a suffix of $q_1 \cdots q_j$ for some $j < z'$.

E.g.)

$$
\begin{array}{ccccccccccccccccc}
1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 & 13 & 14 & 15 & 16 & 17
\end{array}
$$

$$\mathrm{LZ}_{\mathrm{End}}(T) = \mathrm{a}\,|\mathrm{b}|\,\mathrm{a}\;\mathrm{a}\,|\mathrm{b}\;\mathrm{a}\,|\,\mathrm{b}\;\mathrm{b}\,|\,\mathrm{a}\;\mathrm{b}\;\mathrm{b}\;\mathrm{a}\;\mathrm{b}\;\mathrm{a}\;\mathrm{a}\;\mathrm{b}\;\mathrm{b}$$

Suffix of $q_1 q_2$

The longest prefix of $q_5 \cdots$

**Definition:**

The LZ-End factorization of a string $T$ is the factorization $\text{LZ}_{\text{End}}(T) = q_1, \ldots, q_{z'}$ of $T$ such that:

Each phrase $q_i$ $(1 \leq i \leq z' - 1)$ satisfies the following condition.

- $q_i[1, |q_i| - 1]$ is <u>the longest prefix</u> of $q_i \cdots q_{z'}$ which occurs as a suffix of $q_1 \cdots q_j$ for some $j < i$.

The last phrase $q_{z'}$ can be **a suffix** of $T$ which occurs as a suffix of $q_1 \cdots q_j$ for some $j < z'$.

> $z'$ is the number of phrases

E.g.)

$$\text{LZ}_{\text{End}}(T) = \underset{\substack{1\ 2\ 3\ 4\ 5\ 6\ 7\ 8\ 9\ 10\ 11\ 12\ 13\ 14\ 15\ 16\ 17}}{\text{a}|\text{b}|\text{a}\ \text{a}|\text{b}\ \underline{\text{a}|\text{b}}\ \text{b}|\underline{\text{a}\ \text{b}\ \text{b}\ \text{a}}|\text{b}\ \text{a}\ \text{a}\ \text{b}\ \text{b}}$$

Suffix of $q_1 q_2 q_3 q_4 q_5$

The longest prefix of $q_6 \cdots$

19

**Definition:**

The LZ-End factorization of a string $T$ is the factorization $\mathrm{LZ_{End}}(T) = q_1, \ldots, q_{z'}$ of $T$ such that:

$z'$ is the number of phrases

Each phrase $q_i$ ($1 \leq i \leq z' - 1$) satisfies the following condition.

- $q_i[1, |q_i| - 1]$ is <u>the longest prefix</u> of $q_i \cdots q_{z'}$ which occurs as a suffix of $q_1 \cdots q_j$ for some $j < i$.

The last phrase $q_{z'}$ can be **a suffix** of $T$ which occurs as a suffix of $q_1 \cdots q_j$ for some $j < z'$.

E.g.)

$$
\begin{array}{ccccccccccccccccc}
1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 & 13 & 14 & 15 & 16 & 17
\end{array}
$$

$$\mathrm{LZ_{End}}(T) = \text{a}\,|\,\underline{\text{b}}\,|\,\underline{\text{a a}}\,|\,\text{b a}\,|\,\text{b b}\,|\,\text{a b b a}\,|\,\underline{\text{b a a b}}\,|\,\text{b}$$

Suffix of $q_1 q_2 q_3$

The longest prefix of $q_7 \cdots$

**Definition:**

The LZ-End factorization of a string $T$ is the factorization $\mathrm{LZ_{End}}(T) = q_1, \ldots, q_{z'}$ of $T$ such that:

Each phrase $q_i$ $(1 \leq i \leq z' - 1)$ satisfies the following condition.

- $q_i[1, |q_i| - 1]$ is the longest prefix of $q_i \cdots q_{z'}$ which occurs as a suffix of $q_1 \cdots q_j$ for some $j < i$.

The last phrase $q_{z'}$ can be **a suffix** of $T$ which occurs as a suffix of $q_1 \cdots q_j$ for some $j < z'$.

$z'$ is the number of phrases

E.g.)

$$\begin{array}{ccccccccccccccccc} & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 & 13 & 14 & 15 & 16 & 17 \end{array}$$

$$\mathrm{LZ_{End}}(T) = \mathrm{a}\,|\,\underline{\mathrm{b}}\,|\,\mathrm{a}\ \mathrm{a}\,|\,\mathrm{b}\ \mathrm{a}\,|\,\mathrm{b}\ \mathrm{b}\,|\,\mathrm{a}\ \mathrm{b}\ \mathrm{b}\ \mathrm{a}\,|\,\mathrm{b}\ \mathrm{a}\ \mathrm{a}\ \mathrm{b}\,|\,\underline{\mathrm{b}}\,|$$

Suffix of $q_1 q_2$

The longest prefix of $q_8$ and a suffix of $T$

21

**Definition:**

The LZ-End factorization of a string $T$ is the factorization $\mathrm{LZ}_{\mathrm{End}}(T) = q_1, \ldots, q_{z'}$ of $T$ such that:

$z'$ is the number of phrases

Each phrase $q_i$ $(1 \leq i \leq z' - 1)$ satisfies the following condition.

- $q_i[1, |q_i| - 1]$ is <u>the longest prefix</u> of $q_i \cdots q_{z'}$ which occurs as a suffix of $q_1 \cdots q_j$ for some $j < i$.

The last phrase $q_{z'}$ can be **a suffix** of $T$ which occurs as a suffix of $q_1 \cdots q_j$ for some $j < z'$.

E.g.)

$$
\begin{array}{cccccccccccccccccc}
 & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 & 13 & 14 & 15 & 16 & 17 \\
\end{array}
$$

$$\mathrm{LZ}_{\mathrm{End}}(T) = \text{a}|\text{b}|\text{a a}|\text{b a}|\text{b b}|\text{a b b a}|\text{b a a b}|\text{b}|$$

$$q_1 \quad q_2 \quad q_3 \quad\quad q_4 \quad\quad q_5 \quad\quad\quad q_6 \quad\quad\quad\quad q_7 \quad\quad\quad q_8$$

$$\mathrm{z}_{\mathrm{End}}(T) = 8$$

The number of the LZ-End phrases of $T$

22

# The ratio $z_{End} / z_{77}$

It is known that $z_{End}(T) \geq z_{77}(T)$ for any string $T$.
Then how much is the gap between them?
To analyze this, we consider the ratio $z_{End} / z_{77}$.

E.g.)

$$z_{77}(T) = 6$$

$$LZ_{77}(T) = a|b|a\ a|b\ a\ b|b\ a\ b\ b|a\ b\ a\ a\ b\ b|$$

$$LZ_{End}(T) = a|b|a\ a|b\ a|b\ b|a\ b\ b\ a|b\ a\ a\ b|b|$$

$$z_{End}(T) = 8$$

In this case,
$$\frac{z_{End}(T)}{z_{77}(T)} = \frac{8}{6} = 1.333\cdots.$$

# Previous work

**Theorem 1:** [Kreft and Navarro, 2013]

There exist strings $T$ of alphabet size $\sigma = \frac{|T|}{3} + 1$ such that:

$$\frac{z_{\text{End}}(T)}{z_{77}(T)} \to 2 \;\; (|T| \to \infty).$$

# Previous work

**Theorem 1:** [Kreft and Navarro, 2013]

There exist strings $T$ of alphabet size $\sigma = \frac{|T|}{3} + 1$ such that:

$$\frac{z_{\text{End}}(T)}{z_{77}(T)} \to 2 \ (|T| \to \infty).$$

$\Sigma = \{1, 2, \ldots, \sigma\}$

E.g.)
$$T = 1\ 1\ 2 \quad 1\ 1\ 3 \quad 2\ 1\ 4 \quad 3\ 2\ 5 \quad 4\ 3\ 6 \quad \ldots \quad (\sigma - 2)(\sigma - 3)\sigma$$

$$\text{LZ}_{77}(T) = 1\,|\,1\ 2\,|\,1\ 1\ 3\,|\,2\ 1\ 4\,|\,3\ 2\ 5\,|\,4\ 3\ 6\,|\,\ldots\,|\,(\sigma - 2)(\sigma - 3)\sigma\,|$$

$$\text{LZ}_{\text{End}}(T) = 1\,|\,1\ 2\,|\,1\ 1\,|\,3\,|\,2\ 1\,|\,4\,|\,3\ 2\,|\,5\,|\,4\ 3\,|\,6\,|\,\ldots\,|\,(\sigma - 2)(\sigma - 3)\,|\,\sigma\,|$$

**Theorem 1:** [Kreft and Navarro, 2013]

There exist strings $T$ of alphabet size $\sigma = \frac{|T|}{3} + 1$ such that:

$$\frac{\mathrm{z}_{\mathrm{End}}(T)}{\mathrm{z}_{77}(T)} \to 2 \ (|T| \to \infty).$$

$\Sigma = \{1, 2, \dots, \sigma\}$

E.g.)

$$T = 1\ 1\ 2 \quad 1\ 1\ 3 \quad 2\ 1\ 4 \quad 3\ 2\ 5 \quad 4\ 3\ 6 \quad \dots \quad (\sigma - 2)(\sigma - 3)\sigma$$

1

$$\mathrm{LZ}_{77}(T) = 1\,|\,1\ 2\,|\,1\ 1\ 3\,|\,2\ 1\ 4\,|\,3\ 2\ 5\,|\,4\ 3\ 6\,|\,\dots\,|\,(\sigma - 2)(\sigma - 3)\sigma\,|$$

$$\mathrm{LZ}_{\mathrm{End}}(T) = 1\,|\,1\ 2\,|\,1\ 1\,|\,3\,|\,2\ 1\,|\,4\,|\,3\ 2\,|\,5\,|\,4\ 3\,|\,6\,|\,\dots\,|\,(\sigma - 2)(\sigma - 3)\,|\,\sigma\,|$$

2

# Previous work

**Theorem 1:** [Kreft and Navarro, 2013]

There exist strings $T$ of alphabet size $\sigma = \frac{|T|}{3} + 1$ such that:

$$\frac{z_{\text{End}}(T)}{z_{77}(T)} \to 2 \ (|T| \to \infty).$$

$\Sigma = \{1, 2, \ldots, \sigma\}$

E.g.)
$$T = 1\ 1\ 2\ \ \ 1\ 1\ 3\ \ \ 2\ 1\ 4\ \ \ 3\ 2\ 5\ \ \ 4\ 3\ 6\ \ \ \ldots\ \ (\sigma - 2)(\sigma - 3)\sigma$$

$$z_{77}(T) = \sigma$$

$$LZ_{77}(T) = 1\,|\,1\ 2\,|\,1\ 1\ 3\,|\,2\ 1\ 4\,|\,3\ 2\ 5\,|\,4\ 3\ 6\,|\,\ldots\,|\,(\sigma - 2)(\sigma - 3)\sigma\,|$$

$$LZ_{\text{End}}(T) = 1\,|\,1\ 2\,|\,1\ 1\,|\,3\,|\,2\ 1\,|\,4\,|\,3\ 2\,|\,5\,|\,4\ 3\,|\,6\,|\,\ldots\,|\,(\sigma - 2)(\sigma - 3)\,|\,\sigma\,|$$

$$z_{\text{End}}(T) = 2(\sigma - 1)$$

# Previous work

**Theorem 1:** [Kreft and Navarro, 2013]

There exist strings $T$ of alphabet size $\sigma = \frac{|T|}{3} + 1$ such that:

$$\frac{z_{\text{End}}(T)}{z_{77}(T)} \to 2 \ (|T| \to \infty).$$

$\Sigma = \{1, 2, \ldots, \sigma\}$

E.g.)

$$T = 1\,1\,2 \quad 1\,1\,3 \quad 2\,1\,4 \quad 3\,2\,5 \quad 4\,3\,6 \quad \ldots \quad (\sigma - 2)(\sigma - 3)\sigma$$

$$z_{77}(T) = \sigma$$

$$\text{LZ}_{77}(T) = 1\,|\,1\,2\,|\,1\,1\,3\,|\,2\,1\,4\,|\,3\,2\,5\,|\,4\,3\,6\,|\,\ldots\,|\,(\sigma - 2)(\sigma - 3)\sigma\,|$$

$$\text{LZ}_{\text{End}}(T) = 1\,|\,1\,2\,|\,1\,1\,|\,3\,|\,2\,1\,|\,4\,|\,3\,2\,|\,5\,|\,4\,3\,|\,6\,|\,\ldots\,|\,(\sigma - 2)(\sigma - 3)\,|\,\sigma\,|$$

$$z_{\text{End}}(T) = 2(\sigma - 1)$$

$$\frac{z_{\text{End}}(T)}{z_{77}(T)} = \frac{2(\sigma - 1)}{\sigma} = 2 \ (|T| \to \infty, \sigma \to \infty)$$

# Main result

**Theorem 1:** [Kreft and Navarro, 2013]

There exist strings $T$ of alphabet size $\sigma = \frac{|T|}{3} + 1$ such that:
$$\frac{z_{End}(T)}{z_{77}(T)} \to 2 \ (|T| \to \infty).$$

**Theorem 2:** [This work]

There exist strings $S$ of alphabet size $\sigma = \mathbf{2}$ such that:
$$\frac{z_{End}(S)}{z_{77}(S)} \to 2 \ (|S| \to \infty).$$

The string $S$ in Theorem 2 is the period-doubling sequence.

**Definition:**

The $k$-th period-doubling sequence $S_k$ over $\Sigma = \{a, b\}$ is defined as follows:

- $S_0 = a$
- $S_k = S_{k-1} \cdot S_{k-1}[1, n_{k-1}-1] \cdot \overline{c}$   $(k \geq 1)$

$n_{k-1}$ is the length of $S_{k-1}$, that is $n_{k-1} = |S_{k-1}|$.

c is the last character of $S_{k-1}$, that is $c = S_{k-1}[n_{k-1}]$.

$\overline{c}$ is bit-flipped character of c.

Intuition:
Copy the first half and flip the last character.

$S_0 = \text{a}$

Intuition:
Copy the first half and flip the last character.

$S_0 = \boxed{\text{a}}$

$S_1 = \boxed{\text{a}}\,\boxed{\text{b}}$

Intuition:
Copy the first half and flip the last character.

$S_0 = $ a

$S_1 = $ a b

$S_2 = $ a b a a

Intuition:
Copy the first half and flip the last character.

$S_0 = $ a

$S_1 = $ a b

$S_2 = $ a b a a

$S_3 = $ a b a a a b a b

> Intuition:
> Copy the first half and flip the last character.

# Period-doubling sequence [Boston, 1980]

Intuition:
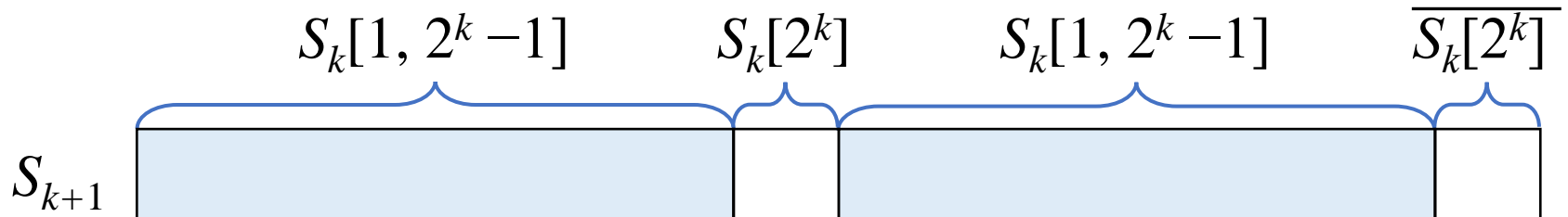Copy the first half and flip the last character.

$S_0 = $ a

$S_1 = $ a b

$S_2 = $ a b a a

$S_3 = $ a b a a a b a b

$S_4 = $ a b a a a b a b a b a a a b a a

⋮

$$S_k[1, 2^k - 1] \quad S_k[2^k] \quad S_k[1, 2^k - 1] \quad \overline{S_k[2^k]}$$

$S_{k+1}$

35

$\text{LZ}_{77}(S_1) = \text{a}|\text{b}|$

$\text{LZ}_{77}(S_2) = \text{a}|\text{b}|\text{a}\,\text{a}|$

$\text{LZ}_{77}(S_3) = \text{a}|\text{b}|\text{a}\,\text{a}|\text{a}\,\text{b}\,\text{a}\,\text{b}|$

$\text{LZ}_{77}(S_4) = \text{a}|\text{b}|\text{a}\,\text{a}|\text{a}\,\text{b}\,\text{a}\,\text{b}|\text{a}\,\text{b}\,\text{a}\,\text{a}\,\text{a}\,\text{b}\,\text{a}\,\text{a}|$

$\text{LZ}_{77}(S_5) = \text{a}|\text{b}|\text{a}\,\text{a}|\text{a}\,\text{b}\,\text{a}\,\text{b}|\text{a}\,\text{b}\,\text{a}\,\text{a}\,\text{a}\,\text{b}\,\text{a}\,\text{a}|\text{a}\,\text{b}\,\text{a}\,\text{a}\,\text{a}\,\text{b}\,\text{a}\,\text{b}\,\text{a}\,\text{b}\,\text{a}\,\text{a}\,\text{a}\,\text{b}\,\text{a}\,\text{b}|$

LZ77 Phrase =

Non-overlapping longest previous occurrence $+$ Single character



$S_k[1, 2^k - 1]$  $S_k[2^k]$  $S_k[1, 2^k - 1]$  $\overline{S_k[2^k]}$

$\text{LZ}_{77}(S_{k+1})$

From definition of period-doubling sequence and LZ77,
$z_{77}(S_k) = k + 1$.

$\text{LZ}_{\text{End}}(S_1) = \text{a}|\text{b}|$

$\text{LZ}_{\text{End}}(S_2) = \text{a}|\text{b}|\text{a a}|$

$\text{LZ}_{\text{End}}(S_3) = \text{a}|\text{b}|\text{a a}|\text{a b a}|\text{b}|$

$\text{LZ}_{\text{End}}(S_4) = \text{a}|\text{b}|\text{a a}|\text{a b a}|\text{b a b}|\text{a a a b a a}|$

$\text{LZ}_{\text{End}}(S_5) = \text{a b a a a b a b a b a a a b a a a b a a a b a b a b a a a b a b}$

# LZ-End of period-doubling sequences $S_k$

$LZ_{End}(S_1) = a|b|$

$LZ_{End}(S_2) = a|b|a\,a|$

$LZ_{End}(S_3) = a|b|a\,a|a\,b\,a|b|$

$LZ_{End}(S_4) = a|b|a\,a|a\,b\,a|b\,a\,b|a\,a\,a\,b\,a\,a|$

$LZ_{End}(S_5) = a|b\,a\,a\,a\,b\,a\,b\,a\,b\,a\,a\,a\,b\,a\,a\,a\,b\,a\,a\,a\,b\,a\,b\,a\,b\,a\,a\,a\,b\,a\,b$

First occurrence

38

$\text{LZ}_{\text{End}}(S_1) = \text{a}|\text{b}|$

$\text{LZ}_{\text{End}}(S_2) = \text{a}|\text{b}|\text{a}\,\text{a}|$

$\text{LZ}_{\text{End}}(S_3) = \text{a}|\text{b}|\text{a}\,\text{a}|\text{a}\,\text{b}\,\text{a}|\text{b}|$

$\text{LZ}_{\text{End}}(S_4) = \text{a}|\text{b}|\text{a}\,\text{a}|\text{a}\,\text{b}\,\text{a}|\text{b}\,\text{a}\,\text{b}|\text{a}\,\text{a}\,\text{a}\,\text{b}\,\text{a}\,\text{a}|$

$\text{LZ}_{\text{End}}(S_5) = \text{a}|\text{b}|\text{a}\,\text{a}\,\text{a}\,\text{b}\,\text{a}\,\text{b}\,\text{a}\,\text{b}\,\text{a}\,\text{a}\,\text{a}\,\text{b}\,\text{a}\,\text{a}\,\text{a}\,\text{b}\,\text{a}\,\text{a}\,\text{a}\,\text{b}\,\text{a}\,\text{b}\,\text{a}\,\text{b}\,\text{a}\,\text{a}\,\text{a}\,\text{b}\,\text{a}\,\text{b}$

First occurrence

# LZ-End of period-doubling sequences $S_k$

$\text{LZ}_{\text{End}}(S_1) = \text{a}|\text{b}|$

$\text{LZ}_{\text{End}}(S_2) = \text{a}|\text{b}|\text{a a}|$

$\text{LZ}_{\text{End}}(S_3) = \text{a}|\text{b}|\text{a a}|\text{a b a}|\text{b}|$

$\text{LZ}_{\text{End}}(S_4) = \text{a}|\text{b}|\text{a a}|\text{a b a}|\text{b a b}|\text{a a a b a a}|$

$\text{LZ}_{\text{End}}(S_5) = \underline{\text{a b}}|\text{a}\,\underline{\text{a}}\,\text{a b a b a b a a a b a a a b a a a b a b a a a b a b}$

Ends with
the 1st phrase

The longest prefix of
the suffix at position 3

# LZ-End of period-doubling sequences $S_k$

$\text{LZ}_{\text{End}}(S_1) = \text{a}|\text{b}|$

$\text{LZ}_{\text{End}}(S_2) = \text{a}|\text{b}|\text{a a}|$

$\text{LZ}_{\text{End}}(S_3) = \text{a}|\text{b}|\text{a a}|\text{a b a}|\text{b}|$

$\text{LZ}_{\text{End}}(S_4) = \text{a}|\text{b}|\text{a a}|\text{a b a}|\text{b a b}|\text{a a a b a a}|$

$\text{LZ}_{\text{End}}(S_5) = \underline{\text{a b}}|\text{a a}|\underline{\text{a b}}|\text{a b a b a a a b a a a b a a a b a b a a a b a b}$

Ends with
the 2nd phrase

The longest prefix of
the suffix at position 5

# LZ-End of period-doubling sequences $S_k$

$\text{LZ}_{\text{End}}(S_1) = \text{a|b|}$

$\text{LZ}_{\text{End}}(S_2) = \text{a|b|a a|}$

$\text{LZ}_{\text{End}}(S_3) = \text{a|b|a a|a b a|b|}$

$\text{LZ}_{\text{End}}(S_4) = \text{a|b|a a|a b a|b a b|a a a b a a|}$

$\text{LZ}_{\text{End}}(S_5) = \text{a|b|a a|a b a|b a b|a a a b a a a b a a a b a b a a a b a b}$

Ends with
the 4th phrase

The longest prefix of
the suffix at position 8

# LZ-End of period-doubling sequences $S_k$

$\mathrm{LZ}_{\mathrm{End}}(S_1) = \text{a|b|}$

$\mathrm{LZ}_{\mathrm{End}}(S_2) = \text{a|b|a a|}$

$\mathrm{LZ}_{\mathrm{End}}(S_3) = \text{a|b|a a|a b a|b|}$

$\mathrm{LZ}_{\mathrm{End}}(S_4) = \text{a|b|a a|a b a|b a b|a a a b a a|}$

$\mathrm{LZ}_{\mathrm{End}}(S_5) = \text{a|b|a a a b a|b a b|a a a b a a|a b a a a b a b a a a b a b}$

Ends with
the 4th phrase

The longest prefix of
the suffix at position 11

# LZ-End of period-doubling sequences $S_k$

$\mathrm{LZ_{End}}(S_1) = \mathrm{a|b|}$

$\mathrm{LZ_{End}}(S_2) = \mathrm{a|b|a\,a|}$

$\mathrm{LZ_{End}}(S_3) = \mathrm{a|b|a\,a|a\,b\,a|b|}$

$\mathrm{LZ_{End}}(S_4) = \mathrm{a|b|a\,a|a\,b\,a|b\,a\,b|a\,a\,a\,b\,a\,a|}$

$\mathrm{LZ_{End}}(S_5) = \underline{\mathrm{a\,b|a\,a|a\,b\,a|b\,a\,b|a\,a\,b\,a\,a}}\underline{\mathrm{a\,b\,a\,a\,b\,a\,b\,a\,b\,a}}\mathrm{|a\,a\,b\,a\,b}$

Ends with
the 5th phrase

The longest prefix of
the suffix at position 17

44

$LZ_{End}(S_1) = a|b|$

$LZ_{End}(S_2) = a|b|a\,a|$

$LZ_{End}(S_3) = a|b|a\,a|a\,b\,a|b|$

$LZ_{End}(S_4) = a|b|a\,a|a\,b\,a|b\,a\,b|a\,a\,a\,b\,a\,a|$

$LZ_{End}(S_5) = a|b|a\,a\,a\,b\,a\,b|a\,b|a\,a\,a\,b\,a\,a|a\,b\,a\,a\,a\,b\,a\,b\,a\,b\,a|a\,a\,b\,a\,b|$

Ends with
the 4th phrase

The longest prefix of
the suffix at position 28

$\text{LZ}_{\text{End}}(S_1) = $ a|b|

$\text{LZ}_{\text{End}}(S_2) = $ a|b|a a|

$\text{LZ}_{\text{End}}(S_3) = $ a|b|a a|a b a|b|

$\text{LZ}_{\text{End}}(S_4) = $ a|b|a a|a b a|b a b|a a a b a a|

$\text{LZ}_{\text{End}}(S_5) = $ a|b|a a|a b a|b a b|a a a b a a|a b a a a b a b a|a a b a b|

$\text{LZ}_{\text{End}}(S_6) = $ a|b|a a|a b a|b a b|a a a b a a|a b a a a b a b a|a a b a b:a b a a|a …

The last phrase of $\text{LZ}_{\text{End}}(S_k)$ is
not always $\text{LZ}_{\text{End}}(S_{k+1})$ phrase.

# LZ-End of period-doubling sequences $S_k$

$\text{LZ}_{\text{End}}(S_1) = \text{a}|\text{b}|$
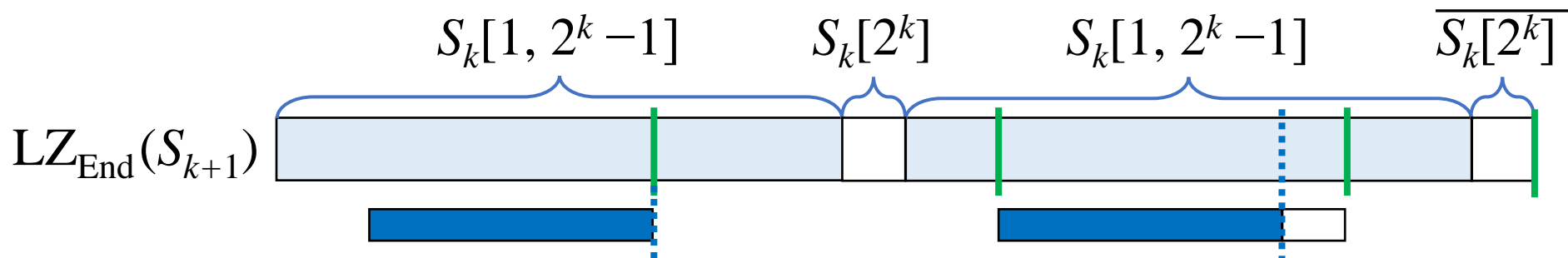
$\text{LZ}_{\text{End}}(S_2) = \text{a}|\text{b}|\text{a}\,\text{a}|$

$\text{LZ}_{\text{End}}(S_3) = \text{a}|\text{b}|\text{a}\,\text{a}|\text{a}\,\text{b}\,\text{a}|\text{b}|$

$\text{LZ}_{\text{End}}(S_4) = \text{a}|\text{b}|\text{a}\,\text{a}|\text{a}\,\text{b}\,\text{a}|\text{b}\,\text{a}\,\text{b}|\text{a}\,\text{a}\,\text{a}\,\text{b}\,\text{a}\,\text{a}|$

$\text{LZ}_{\text{End}}(S_5) = \text{a}|\text{b}|\text{a}\,\text{a}|\text{a}\,\text{b}\,\text{a}|\text{b}\,\text{a}\,\text{b}|\text{a}\,\text{a}\,\text{a}\,\text{b}\,\text{a}\,\text{a}|\text{a}\,\text{b}\,\text{a}\,\text{a}\,\text{a}\,\text{b}\,\text{a}\,\text{b}\,\text{a}|\text{a}\,\text{b}\,\text{a}\,\text{b}|$

$\text{LZ}_{\text{End}}(S_6) = \text{a}|\text{b}|\text{a}\,\text{a}|\text{a}\,\text{b}\,\text{a}|\text{b}\,\text{a}\,\text{b}|\text{a}\,\text{a}\,\text{a}\,\text{b}\,\text{a}\,\text{a}|\text{a}\,\text{b}\,\text{a}\,\text{a}\,\text{a}\,\text{b}\,\text{a}\,\text{b}\,\text{a}|\text{a}\,\text{b}\,\text{a}\,\text{b}|\text{a}\,\text{b}\,\text{a}\,\text{a}\ldots$



$\text{LZ}_{\text{End}}(S_{k+1})$

$S_k[1, 2^k - 1] \qquad S_k[2^k] \qquad S_k[1, 2^k - 1] \qquad \overline{S_k[2^k]}$

From definition of period-doubling sequence and LZ-End,
$z_{\text{End}}(S_k) = 2k - \text{O}(\log^* k).$

**Observation 1:**

$LZ_{End}(S_4)$

a|b|a|a|ab|ab|ab|aaabaa|

$LZ_{End}(S_5)$

a|b|a|a|ab|ab|ab|aaabaa|abaaababaa|abab|

$LZ_{End}(S_6)$

a|b|a|a|ab|ab|ab|aaabaa|abaaabababa|aababaabaa|ababaaabaaabaaabababaa|abaa|

**Observation 1:**

$LZ_{End}(S_4)$

a|b|a|a|b|a|b|ab|aaabaa|

$LZ_{End}(S_5)$

a|b|a|a|b|a|b|ab|aaabaa|abaaababa|aabab|

$LZ_{End}(S_6)$

a|b|a|a|b|a|b|ab|aaabaa|abaaababa|aabababaa|ababaaabaaabaaababaa|abaa|

**Observation 1:**

$\mathrm{LZ}_{\mathrm{End}}(S_4)$

a|b|a|a|b|a|b|ab|aaabaa|

$\qquad\qquad\qquad 6$

$\mathrm{LZ}_{\mathrm{End}}(S_5)$

a|b|a|a|b|a|b|ab|aaabaa|abaaababa|aabab|

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad 5$

$\mathrm{LZ}_{\mathrm{End}}(S_6)$

a|b|a|a|b|a|b|ab|aaabaa|abaaababa|aababaaa|babababaa|abababaaabaaabaaabababaa|abaa|

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad 4$

The length of the last LZ-End phrase decreases by **1**
until the length of the last phrase becomes **1**.

**Observation 2:**

$S_4$

```
   LZ77:abaaababaaabaaabaa
LZ-End:abaaababababaaabaa
```

$S_5$

```
   LZ77:abaaababaaaabaaabaaababababaaabab
LZ-End:abaaababababaaabaaabaaababababaaabab
```

$S_6$

```
   LZ77:abaaababaaaabaaabaaababababaaababaaaababababaaabaaabaaababababaaabaa
LZ-End:abaaababababaaabaaabaaababababaaababaaabababaaabababaaabaaabaaababababaaabaa
```

**Observation 2:**

$S_4$  $z_{77}(S_4) = 5, z_{End}(S_4) = 6$

LZ77: a b a a a b a b a b a a a b a a

LZ-End: a b a a a b a b a b a a a b a a

$S_5$  $z_{77}(S_5) = 6, z_{End}(S_5) = 8$

LZ77: a b a a a b a b a b a a a b a a a b a a a b a b a b a a a b a b

LZ-End: a b a a a b a b a b a a a b a a a b a a a b a b a b a a a b a b

$S_6$  $z_{77}(S_6) = 7, z_{End}(S_6) = 10$

LZ77: a b a a a b a b a b a a a b a a a b a a a b a b a b a a a b a b a b a a a b a b a b a a a b a a a b a a a b a b a b a a a b a a

LZ-End: a b a a a b a b a b a a a b a a a b a a a b a b a b a a a b a b a b a a a b a b a b a a a b a a a b a a a b a b a b a a a b a a

# Increase of LZ-End phrases

**Observation 2:**

Increasing number of phrases:
- LZ77: **1** (for any $k$)
- LZ-End: **2** (for almost all $k$)

$S_4$  $z_{77}(S_4) = \mathbf{5}$, $z_{\text{End}}(S_4) = \mathbf{6}$

LZ77: abaaababaaabaa

LZ-End: abaaababaaabaa

LZ77: + **1**
LZ-End: + **2**

$S_5$  $z_{77}(S_5) = \mathbf{6}$, $z_{\text{End}}(S_5) = \mathbf{8}$

LZ77: abaaababaaabaaabaaababaaabab

LZ-End: abaaababaaabaaabaaababaaabab

LZ77: + **1**
LZ-End: + **2**

$S_6$  $z_{77}(S_6) = \mathbf{7}$, $z_{\text{End}}(S_6) = \mathbf{10}$

LZ77: abaaababaaabaaabaaababaaababaaabababaaabaaabaaabababaaabaa

LZ-End: abaaababaaabaaabaaabababaaababaaabababaaabaaabaaabababaaabaa

# Table of LZ77 and LZ-End phrases

$$z_{End}(S_k) - z_{End}(S_{k-1})$$

| $k$ | $z_{End}$ | $z_{77}$ | $z_{End}$ / $z_{77}$ | Length of the last LZ-End phrase | $z_{End}$ diff |
|---|---|---|---|---|---|
| ... | ... | ... | ... | ... | ... |
| 6 | 10 | 7 | 1.428... | 4 | 2 |
| 7 | 12 | 8 | 1.5 | 3 | 2 |
| 8 | 14 | 9 | 1.555... | 2 | 2 |
| 9 | 16 | 10 | 1.6 | 1 | 2 |
| 10 | 17 | 11 | 1.545... | 384 | 1 |
| 11 | 19 | 12 | 1.583... | 383 | 2 |
| ... | ... | ... | ... | ... | ... |
| 393 | 783 | 394 | 1.987... | 1 | 2 |
| 394 | 784 | 395 | 1.984... | $3*2^{391}$ | 1 |
| 395 | 786 | 396 | 1.984... | $3*2^{391} - 1$ | 2 |
| ... | ... | ... | ... | ... | ... |
| $3*2^{391} + 394$ | ... | ... | 1.999... | $3*2^{(3*2^{391} + 391)}$ | 1 |

# Table of LZ77 and LZ-End phrases

$z_{End}(S_k) - z_{End}(S_{k-1})$

| $k$ | $z_{End}$ | $z_{77}$ | $z_{End} / z_{77}$ | Length of the last LZ-End phrase | $z_{End}$ diff |
|---|---|---|---|---|---|
| ... | ... | ... | ... | ... | ... |
| 6 | 10 | 7 | 1.428... | 4 | 2 |
| 7 | 12 | 8 | 1.5 | 3 | 2 |
| 8 | 14 | 9 | 1.555... | 2 | 2 |
| 9 | 16 | 10 | 1.6 | 1 | 2 |
| 10 | 17 | 11 | 1.545... | 384 | 1 |
| 11 | 19 | 12 | 1.583... | 383 | 2 |
| ... | ... | ... | ... | ... | ... |
| 393 | 783 | 394 | 1.987... | 1 | 2 |
| 394 | 784 | 395 | 1.984... | $3*2^{391}$ | 1 |
| 395 | 786 | 396 | 1.984... | $3*2^{391} - 1$ | 2 |
| ... | ... | ... | ... | ... | ... |
| $3*2^{391} + 394$ | ... | ... | 1.999... | $3*2^{(3*2^{391} + 391)}$ | 1 |

# Table of LZ77 and LZ-End phrases

$$z_{\text{End}}(S_k) - z_{\text{End}}(S_{k-1})$$

| $k$ | $z_{\text{End}}$ | $z_{77}$ | $z_{\text{End}} / z_{77}$ | Length of the last LZ-End phrase | $z_{\text{End}}$ diff |
|---|---|---|---|---|---|
| ... | ... | ... | ... | ... | ... |
| 6 | 10 | 7 | 1.428... | 4 | 2 |
| 7 | 12 | 8 | 1.5 | 3 | 2 |
| 8 | 14 | 9 | 1.555... | 2 | 2 |
| 9 | 16 | 10 | 1.6 | 1 | 2 |
| 10 | 17 | 11 | 1.545... | 384 | 1 |
| 11 | 19 | 12 | 1.583... | 383 | 2 |
| ... | ... | ... | ... | ... | ... |
| 393 | 783 | 394 | 1.987... | 1 | 2 |
| 394 | 784 | 395 | 1.984... | $3*2^{391}$ | 1 |
| 395 | 786 | 396 | 1.984... | $3*2^{391} - 1$ | 2 |
| ... | ... | ... | ... | ... | ... |
| $3*2^{391} + 394$ | ... | ... | 1.999... | $3*2^{(3*2^{391} + 391)}$ | 1 |

# Table of LZ77 and LZ-End phrases

$$z_{End}(S_k) - z_{End}(S_{k-1})$$

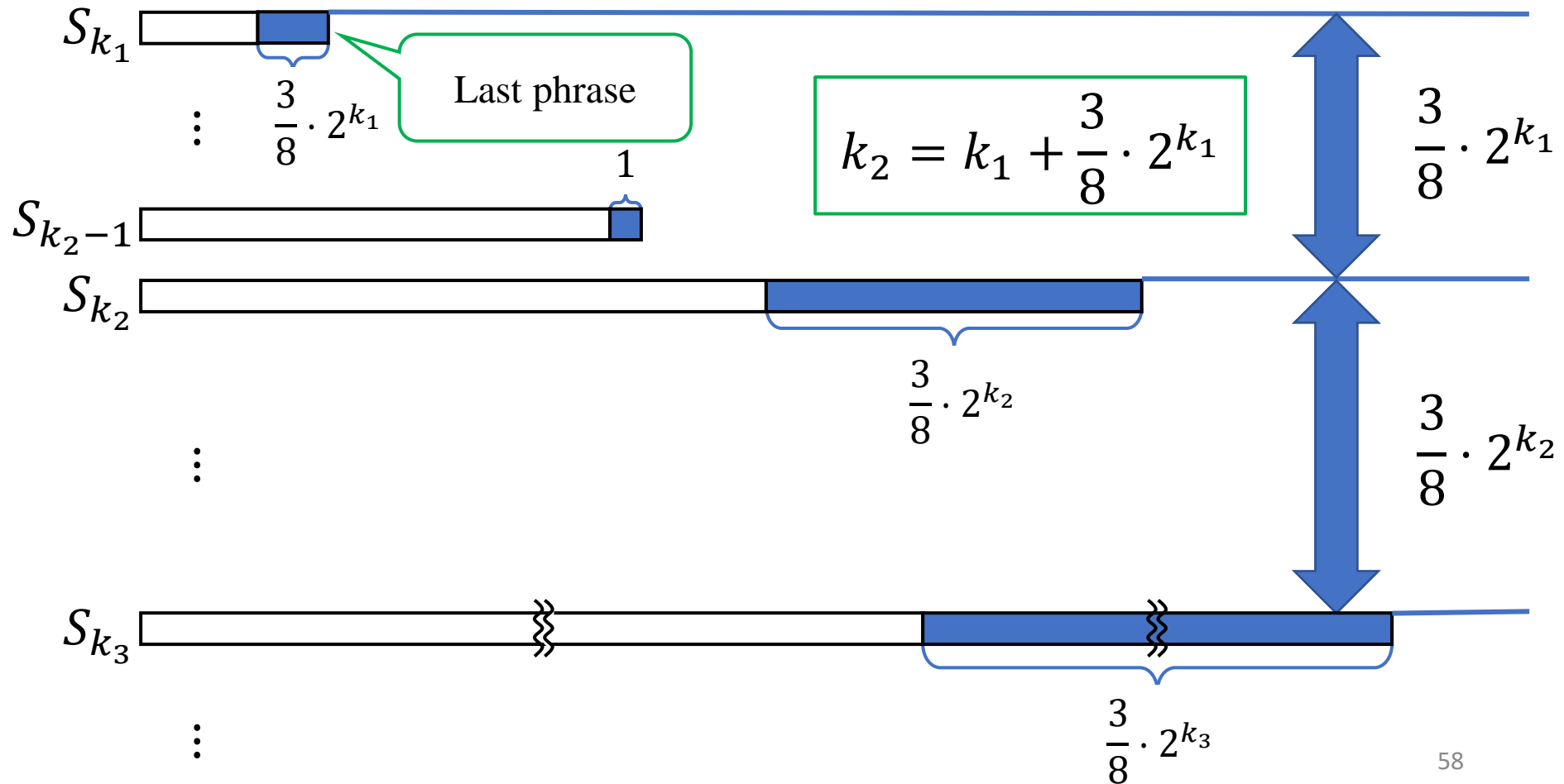| $k$ | $z_{End}$ | $z_{77}$ | $z_{End} / z_{77}$ | Length of the last LZ-End phrase | $z_{End}$ diff |
|---|---|---|---|---|---|
| ... | ... | ... | ... | ... | ... |
| 6 | 10 | 7 | 1.428... | 4 | 2 |
| 7 | 1? | | | | 2 |
| 8 | 1 | | | | 2 |
| 9 | 1 | | | | 2 |
| 10 | 1 | | | | 1 |
| 11 | 19 | 12 | 1.583... | 383 | 2 |
| ... | ... | ... | ... | ... | ... |
| 393 | 783 | 394 | 1.987... | 1 | 2 |
| 394 | 784 | 395 | 1.984... | $3*2^{391}$ | 1 |
| 395 | 786 | 396 | 1.984... | $3*2^{391} - 1$ | 2 |
| ... | ... | ... | ... | ... | ... |
| $3*2^{391} + 394$ | ... | ... | 1.999... | $3*2^{(3*2^{391} + 391)}$ | 1 |

**Lemma 1:**
The number of **1**'s is $O(\log^* k)$.
Thus, $z_{End}(S_k) = 2k - O(\log^* k)$.

57

# Why $O(\log^* k)$?

$$k_m = O\left(2^{2^{\cdot^{\cdot^{2^k}}}}\right) \Leftrightarrow m = O(\log^* k)$$

**Lemma 2:**

The maximal length of the last LZ-End phrase is $\frac{3}{8} \cdot 2^k$.



$S_{k_1}$

$\frac{3}{8} \cdot 2^{k_1}$

Last phrase

1

$k_2 = k_1 + \frac{3}{8} \cdot 2^{k_1}$

$\frac{3}{8} \cdot 2^{k_1}$

$S_{k_2-1}$

$S_{k_2}$

$\frac{3}{8} \cdot 2^{k_2}$

$\frac{3}{8} \cdot 2^{k_2}$

$S_{k_3}$

$\frac{3}{8} \cdot 2^{k_3}$

58

# The ratio $z_{End} / z_{77}$

We obtain the following result.

$$z_{77}(S_k) = k + 1$$

$$z_{End}(S_k) = 2k - O(\log^* k)$$

$$\frac{z_{End}(S_k)}{z_{77}(S_k)} = \frac{2k - O(\log^* k)}{k + 1} \to 2 \quad (k \to \infty).$$

Period-doubling sequence

**Theorem 2:**

There exist strings $S$ of alphabet size $\sigma = $ **2** such that:
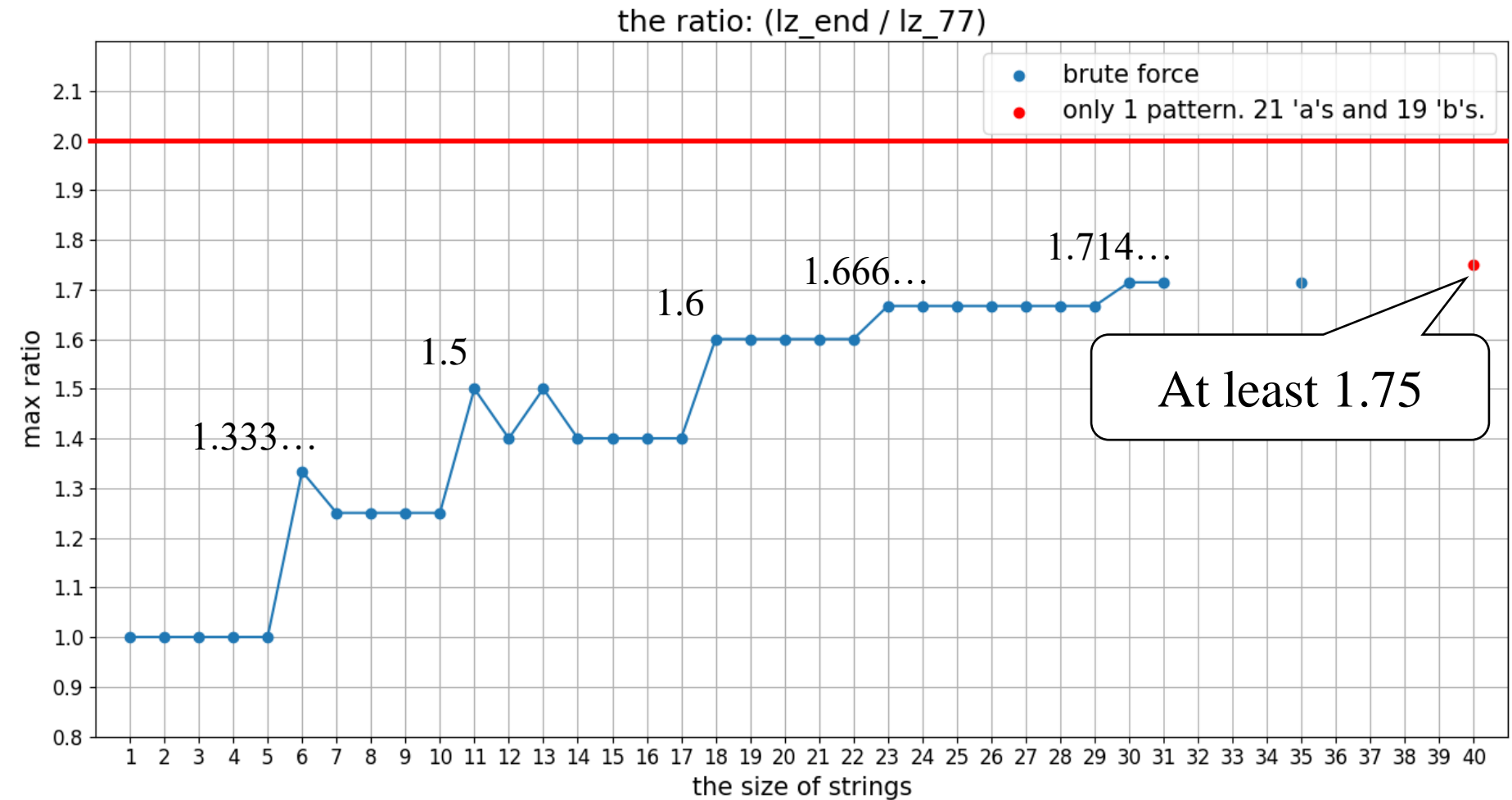$$\frac{z_{End}(S)}{z_{77}(S)} \to 2 \ (|S| \to \infty).$$

# Summary and future work

Summary:

- We proved that period-doubling sequence $S$ satisfies that $z_{\text{End}}(S) / z_{77}(S)$ asymptotically approaches 2 when the limit as the length of $S$ tends to infinity.
- There also exist other binary sequences $S'$ such that $z_{\text{End}}(S') / z_{77}(S')$ asymptotically approaches 2.

**Conjecture:** [Kreft and Navarro, 2013]

$z_{\text{End}}(T) / z_{77}(T) \leq 2$ holds for any string $T$.

the ratio: (lz_end / lz_77)

The ratio seems to asymptotically approach 2.

# Summary and future work

Summary:

- We proved that period-doubling sequence $S$ satisfies that $z_{\text{End}}(S) / z_{77}(S)$ asymptotically approaches 2 when the limit as the length of $S$ tends to infinity.
- There also exist other binary sequences $S'$ such that $z_{\text{End}}(S') / z_{77}(S')$ asymptotically approaches 2.

Future work:

- Prove or disprove the conjecture for upper bound.

**Conjecture:** [Kreft and Navarro, 2013]

$z_{\text{End}}(T) / z_{77}(T) \leq 2$ holds for any string $T$.