

RePair Grammars are the Smallest Grammars for Fibonacci Words

Takuya Mieno (University of Electro-Communications)

Shunsuke Inenaga (Kyushu University)

Takashi Horiyama (Hokkaido University)



Straight Line Program (SLP)

A context-free grammar in the Chomsky normal form that produces only a single string w is called a **straight-line program (SLP)** for w .

E.g.

The grammar $G = (V, \Sigma, \delta, S)$ below is an SLP which produces a single string $w = ababaabaaba$.

$$V = \{X_7, X_6, X_5, X_4, X_3, X_2, X_1, A, B\}$$

$$\Sigma = \{a, b\}, S = X_7,$$

$$\delta = \{$$

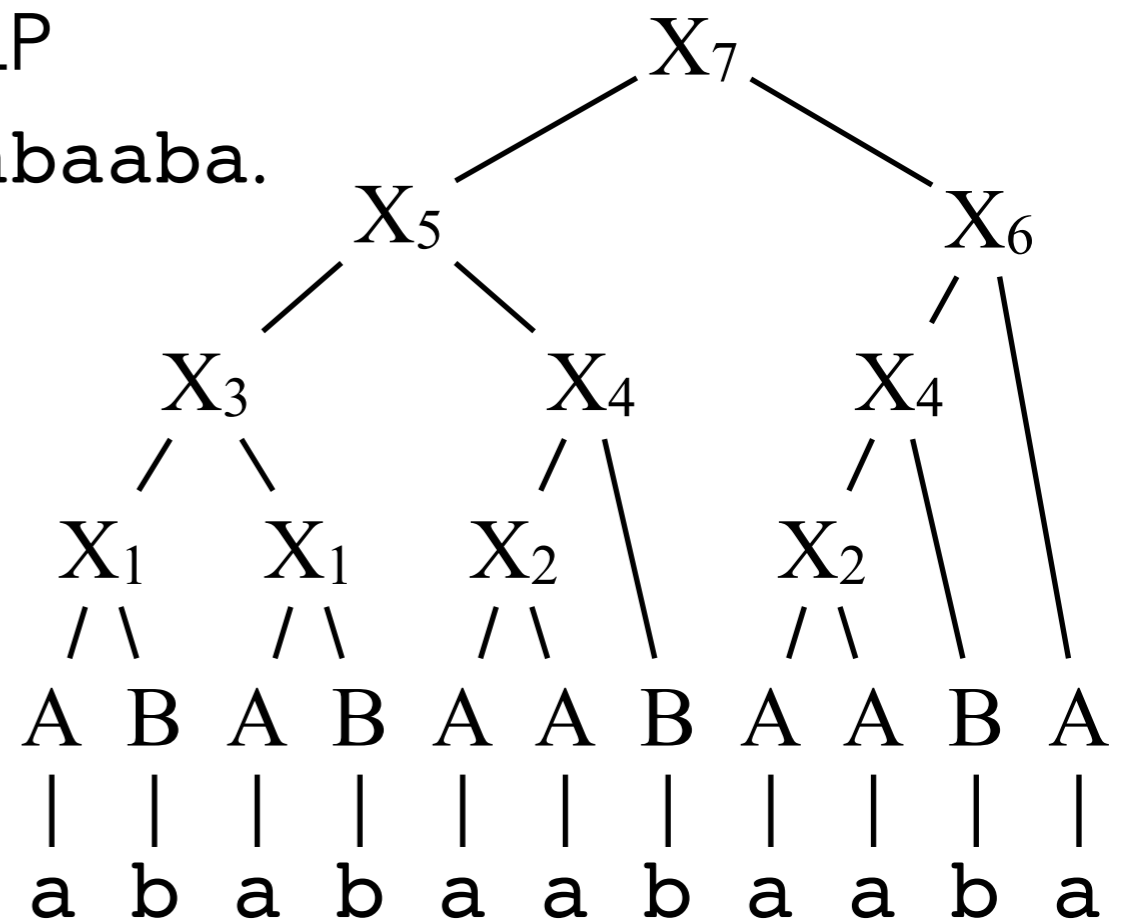
$$X_7 \rightarrow X_5X_6, X_6 \rightarrow X_4A, X_5 \rightarrow X_3X_4,$$

$$X_4 \rightarrow X_2B, X_3 \rightarrow X_1X_1, X_2 \rightarrow AA,$$

$$X_1 \rightarrow AB, A \rightarrow a, B = b$$

$$\}$$

$$|G| = 9$$



The **size** $|G|$ of an SLP G is **# of productions** in G .

Smaller grammar is better in terms of string compression.

Smallest Grammar Problem (SLP version)

Input : String w

Output : SLP for w with the **smallest** size.

This problem is known to be NP-hard [Charikar et al., '05].

Even if the alphabet size σ is a constant with $\sigma \geq 17$ [Casel et al., '21].

The hardness for $1 \leq \sigma \leq 16$ is open.

There are practical grammar compressors.

RePair, LZ78, BISECTION, SEQUENTIAL, etc.

Bounds for the *approximation ratios* of these compressors have been investigated [Charikar et al., '05], [Bannai et al., '21].

Our Result (informal statement)

We focus on the smallest grammars of **Fibonacci words**:

$$F_1 = \mathbf{b}, F_2 = \mathbf{a}, F_i = F_{i-1}F_{i-2} \ (i \geq 3).$$

n	F_n	length
1	b	1
2	a	1
3	ab	2
4	aba	3
5	abaab	5
6	abaababa	8

Our main result

For any **Fibonacci word** F_n , the set of **smallest grammars** of F_n **equals** the set of grammars obtained by applying **RePair** algorithms to F_n .

To our knowledge, Fibonacci words are **the first non-trivial strings** whose smallest grammar sizes are computable in polynomial time (except trivial ones such as a^{2^k} and $(ab)^{2^k}$).

RePair [Larsson & Moffat, '99]

RePair is a greedy algorithm which computes an SLP for given string w .

1. Replace all terminals in w with non-terminals.
2. Recursively do the following while the length of string is ≥ 2 :
Choose a **most frequent bigram** and replace it with a non-terminal.

E.g.

$w = ababaaaaab$

Replace terminals a, b with non-terminals A, B.

ABABAAAAAB

of occurrences of each bigram are

AA: 2, AB: **3**, BA: 2, BB: 0.

→ choose AB and replace it with non-terminal C.

$\delta = \{$
 $A \rightarrow a,$ CCAAAAAC
 $B \rightarrow b,$
 $C \rightarrow AB,$ CCDDC
 $D \rightarrow AA,$
 $E \rightarrow CC,$ EDDC
 \vdots
 $\}$

of occurrences of all bigrams are the same, so we can choose a most frequent one **arbitrarily**. Here, choose CC and replace it with E.

RePair Grammars

RePair is a greedy algorithm which computes an SLP for given string w . When there are multiple frequent bigram, we can choose any of them **arbitrarily**.

Which one is chosen depends on the implementation of RePair.

E.g.

Most frequent bigrams in $w = \text{ababaabaaba}$ are **ab and ba**.

If we replace **ab** with X , the string changes to XXaXaXa .

If we replace **ba** with X , the string changes to aXXaXaX .

A grammar which is obtained by applying (any implementation of) RePair to w is called a **RePair grammar** of w .

We denote by $\text{RePair}(w)$ the set of all RePair grammars of w .

Our Result (formal statement)

Let $\text{Opt}(w)$ be the set of smallest grammars of w .

Theorem 1 [This work]

For every $n \geq 1$, $\text{Opt}(F_n) = \text{RePair}(F_n)$ holds. Also, for every $n \geq 4$, $|\text{Opt}(F_n)| = n - 2$ if n is even, and $|\text{Opt}(F_n)| = n - 1$ if n is odd.

The proof consists of 3 steps. We will show that

Lemma 1



the smallest size of grammars of F_n equals n ,

Lemma 2

the size of **any** RePair grammar of F_n is the smallest, and

Lemma 3

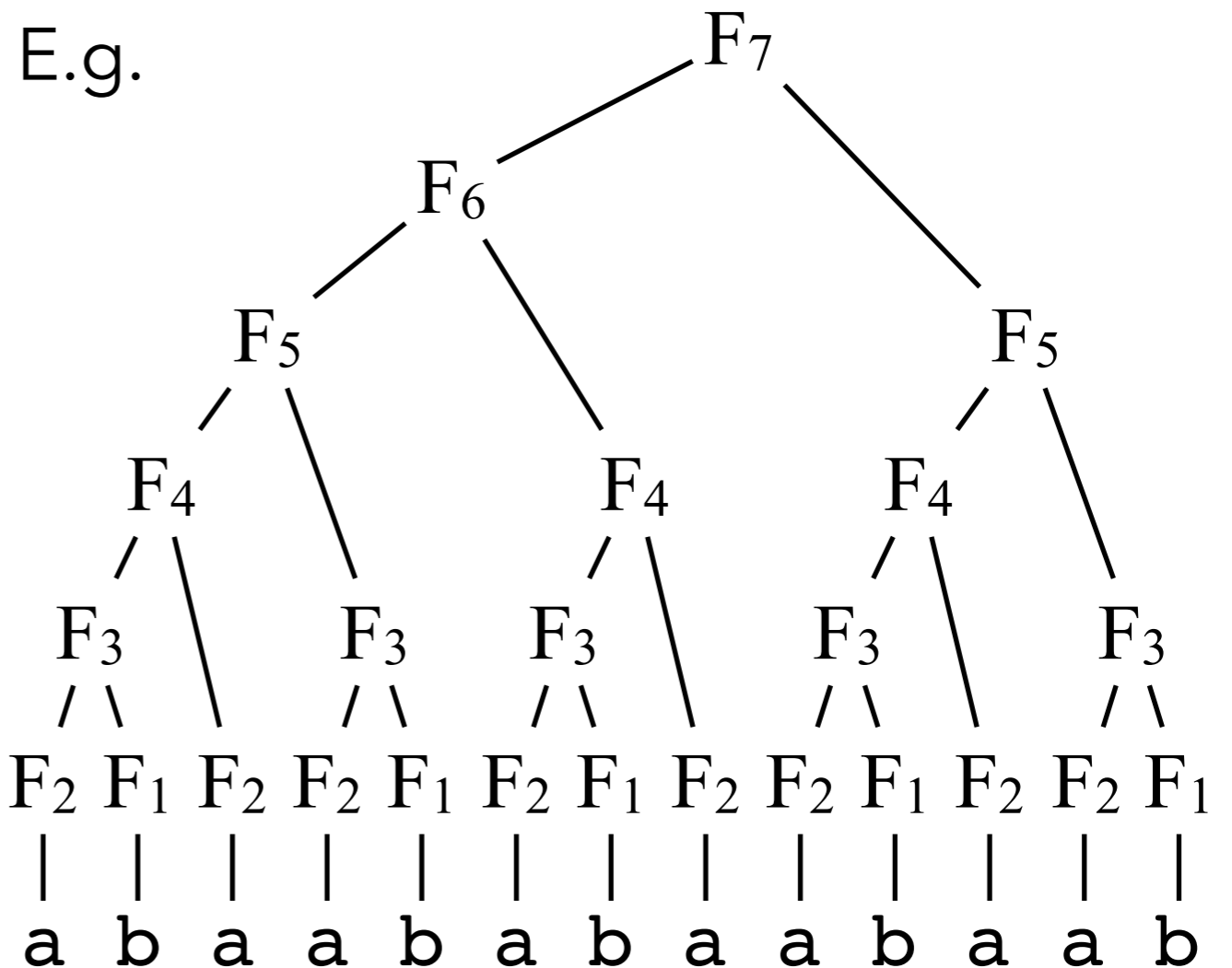
the size of **any non-RePair** grammar of F_n is **not** the smallest.

Smallest Grammar Size of Fibonacci Words

There is a size- n grammar by the definition of F_n .

$$F_1 = \mathbf{b}, F_2 = \mathbf{a},$$

$$F_i = F_{i-1}F_{i-2} \ (i \geq 3).$$



The remaining task is to prove that there is no grammar of size $< n$.

Relation Between Grammar and LZ-factorization

To evaluate the grammar size of F_n , we utilize the **LZ-factorization**.

Definition (LZ-factorization (without self-reference))

The **LZ-factorization** of string w is a sequence (p_1, \dots, p_z) of strings such that $w = p_1 \dots p_z$ and each p_i is either a fresh symbol or the longest prefix of $p_i \dots p_z$ which occurs in $p_1 \dots p_{i-1}$.

Let $z(w)$ be the number of factors in the LZ-factorization of w .

E.g.

$w = a|b|a|a\ b\ a|a\ b|c|a\ b\ a\ a|c\ a\ b\ a\ a|b\ a$

The size of LZ-factorization is $z(w) = 9$

Theorem 2 [Rytter, '03]

For any string w , $z(w) \leq g^*(w)$ holds.

$g^*(w)$: the size of the smallest grammar of w .

Proof of Lemma 1

Theorem 2 [Rytter, '03]

For any string w , $z(w) \leq g^*(w)$ holds.

With a modification to the proof of the Theorem 2, we obtain a slightly tighter lower bound on the smallest grammar size:

Theorem 3 [This work]

For any string w , $z(w) - 1 + \sigma_w \leq g^*(w)$ holds
where σ_w is the number of characters appearing in w .

Since $z(F_n) = n - 1$ holds for Fibonacci words,

$z(F_n) - 1 + \sigma_w = n - 1 - 1 + 2 = n \leq g^*(F_n)$ holds by Theorem 3.

Namely, the aforementioned grammar size n is the smallest. \square

Lemma 1

The size of the smallest grammar of F_n is n .

Lemmas for Our Main Result

Lemma 1

The size of the smallest grammar of F_n is n .

Lemma 2 ◀ NEXT

The size of **any** RePair grammar of F_n is n .

By Lemmas 1 and 2, $\text{RePair}(F_n) \subseteq \text{Opt}(F_n)$ holds.

Lemma 3

The size of **any non-RePair grammar** of F_n is at least $n + 1$.

By combining Lemma 3, we obtain $\text{RePair}(F_n) = \text{Opt}(F_n)$.

Most Frequent Bigrams in Fibonacci Words

The most frequent bigrams in F_n are ab and/or ba.

n	F_n	most frequent bigram(s)
3	ab	ab
4	aba	ab, ba
5	abaab	ab
6	abaababa	ab, ba
7	abaababaabaab	ab
8	abaababaabaababa	ab, ba
9	abaababaabaababaabaabaabaabaab	ab
10	abaababaabaababaabaabaabaabaabaabaabaabaabaabaabaabaabaabaab	ab, ba

E.g. # of occurrences of bigrams in $F_7 = \text{abaababaabaab}$
 ab: 5, ba: 4, aa: 3, bb: 0

There are two possible cases, depending on which bigram is chosen.

When ab is Chosen

Let's observe how Fibonacci words change when ab is replaced to another symbol.

E.g. The most frequent bigram in F_7 is (only) ab .

$$F_7 = \underline{ab} \underline{a} \underline{ab} \underline{ab} \underline{a} \underline{ab} \underline{ab}$$

Replace all ab 's with symbol X , then the string changes to

$$\underline{X} \underline{a} \underline{X} \underline{X} \underline{a} \underline{X} \underline{a} \underline{X}$$

$$F_6 = \underline{a} \underline{b} \underline{a} \underline{a} \underline{b} \underline{a} \underline{b} \underline{a}$$

This is *isomorphic* to $F_6 = abaababa$.

In general, Lemma 4 holds.

Lemma 4

By replacing all occurrences of ab in F_n with a new symbol, we obtain a string which is isomorphic to F_{n-1} .

When ab is Chosen

Lemma 4

By replacing all occurrences of ab in F_n with a new symbol, we obtain a string which is isomorphic to F_{n-1} .

Lemma 4 can be easily proven by using another definition of F_n :

$$F_n = \phi^{n-1}(\mathbf{b})$$

where ϕ is the *string morphism* such that $\phi(\mathbf{a}) = \mathbf{ab}$ and $\phi(\mathbf{b}) = \mathbf{a}$.

E.g.

$$F_6 = \underline{\mathbf{a}} \mathbf{b} \underline{\mathbf{a}} \underline{\mathbf{a}} \mathbf{b} \underline{\mathbf{a}} \mathbf{b} \underline{\mathbf{a}}$$

$\phi \downarrow$

\uparrow Replace $ab \rightarrow a$
(and then $a \rightarrow b$)

$$F_7 = \underline{\mathbf{ab}} \underline{\mathbf{a}} \underline{\mathbf{abab}} \underline{\mathbf{a}} \underline{\mathbf{ab}} \underline{\mathbf{a}} \underline{\mathbf{ab}}$$

Essentially, the replacement in Lemma 4 is the **inverse of ϕ** .

→ Of course, the resulting string is isomorphic to F_{n-1} . □

When ba is Chosen

Let's observe how Fibonacci words change when ba is replaced with another symbol.

E.g. The most frequent bigrams in F_8 are ab and ba .

$F_8 = \text{abaababaabaababaababa}$

Replace all ba 's with symbol X , then the string changes to

$F_7 = \text{a X a X X a X a X X a X X}$
 $\text{a b a a b a b a a b a a b}$

This is isomorphic to the *right-rotation* R_7 of F_7 .

$$R_n := F_n[|F_n|] \cdot F_n[1.. |F_n|-1]$$

In general, Lemma 5 holds.

Lemma 5

By replacing all occurrences of ba in F_{2k} with a new symbol, we obtain a string which is isomorphic to the right-rotation R_{2k-1} of F_{2k-1} .

Applying RePair to Right-rotation of Fibonacci Words (1/2)

n	$R_n := F_n[F_n] \cdot F_n[1.. F_n -1]$	most frequent bigram
3	ba	ba
4	aab	ab
5	babaa	ba
6	aabaabab	ab
7	babaababaabaa	ba
8	aabaababaabaababab	ab
9	babaababaabaababaabaababaabaa	ba
10	aabaababaabaababaabaabaabaabaabaabaabaabaabaabaabaabab	ab

$R_7 = \text{babaababaabaa}$
 $\rightarrow \text{XXaXXaXa}$
 (aabaabab)

Lemma 6

By replacing all occurrences of ba in R_{2k+1} with a new symbol, we obtain a string which is isomorphic to R_{2k} .

Applying RePair to Right-rotation of Fibonacci Words (2/2)

n	$R_n := F_n[F_n] \cdot F_n[1.. F_n -1]$	most frequent bigram
3	ba	ba
4	aab	ab
5	babaa	ba
6	aabaabab	ab
7	babaababaabaa	ba
8	aabaababaabaabaabab	ab
9	babaababaabaabaabaabaabaabaa	ba
10	aabaababaabaabaabaabaabaabaabaabaabaabaabaabaabaabab	ab

$R_6 = \text{aabaabab}$
 $\rightarrow \text{aXaXX}$
 (babaa)

Lemma 7

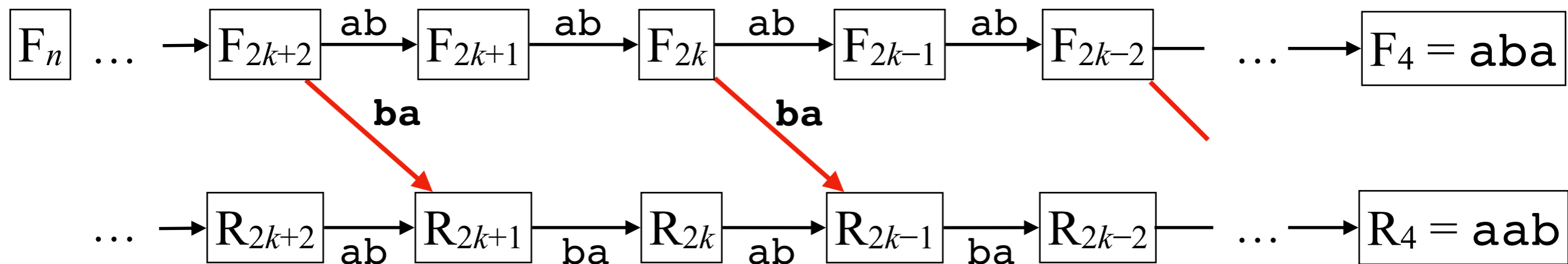
By replacing all occurrences of ab in R_{2k} with a new symbol, we obtain a string which is isomorphic to R_{2k-1} .

"RePair Graph" for Fibonacci words (1/2)

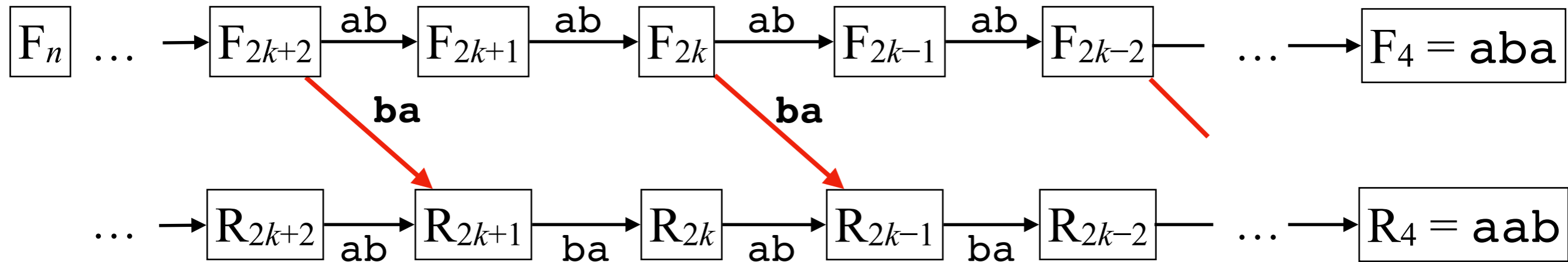
From Lemmas 4–7, the strings that appear during the execution of RePair for a Fibonacci word are strings that are isomorphic to

- some Fibonacci word F_i , or
- the right-rotation R_i of some Fibonacci word.

The **changes** of strings when RePair is applied to Fibonacci words can be represented by the following "**RePair graph**".



"RePair Graph" for Fibonacci words (2/2)



The size of a grammar equals **the length of the corresponding path** from the left-end (F_n) to the right-end (F_4 or R_4) of the RePair graph.

→ That always equals n , i.e., the smallest size.

The number of RePair grammars equals **the number of paths** from the left-end (F_n) to the right-end (F_4 or R_4) of the RePair graph.

→ That equals $n - 2$ if n is even, or $n - 1$ otherwise. \square

Lemmas for Our Main Result

Lemma 1 ✓ (omit the details in this talk)

The size of the smallest grammar of F_n is n .

Lemma 2 ✓

The size of **any** RePair grammar of F_n is n .

By Lemmas 1 and 2, $\text{RePair}(F_n) \subseteq \text{Opt}(F_n)$ holds.

Lemma 3 ◀ **NEXT**

The size of **any non-RePair grammar** of F_n is at least $n + 1$.

By combining Lemma 3, we obtain $\text{RePair}(F_n) = \text{Opt}(F_n)$.

16 Strategies in non-RePair Algorithms

Bigram to replace Target string		aa		ab		ba	
		all	not all	all	not all	all	not all
F_{2k}	1	2	RePair	3	RePair	4	
F_{2k+1}			RePair		5	6	
R_{2k}	7	8	RePair	9	10	11	
R_{2k+1}	12	13	14	15	RePair	16	

For all **16 strategies** to replace bigrams which do not satisfy the RePair conditions, the resulting grammars are non-smallest.

To prove this, we heavily utilize the fact "The smallest grammar size is lower-bounded by the LZ77 size" (Theorem 3).

Lemmas for Our Main Result

Lemma 1 ✓ (omit the details in this talk)

The size of the smallest grammar of F_n is n .

Lemma 2 ✓

The size of **any** RePair grammar of F_n is n .

By Lemmas 1 and 2, $\text{RePair}(F_n) \subseteq \text{Opt}(F_n)$ holds.

Lemma 3 ✓

The size of **any non-RePair grammar** of F_n is at least $n + 1$.

By combining Lemma 3, we obtain $\text{RePair}(F_n) = \text{Opt}(F_n)$.

Conclusion and Future Work

Conclusion

We showed that the **smallest grammars** are the same as the **RePair grammars** of Fibonacci words F_n .

The smallest grammars of F_n are characterized completely.

Theorem 1 [This work]

For every $n \geq 1$, $\text{Opt}(F_n) = \text{RePair}(F_n)$ holds. Also, for every $n \geq 4$, $|\text{Opt}(F_n)| = n - 2$ if n is even, and $|\text{Opt}(F_n)| = n - 1$ if n is odd.

Future Work

To investigate whether it is possible to characterize the smallest grammars of other binary words, including Thue-Morse words TM_n and Period-doubling words PD_n . We have confirmed that

$$\text{Opt}(\text{TM}_n) \neq \text{RePair}(\text{TM}_n) \text{ and } \text{Opt}(\text{PD}_n) \neq \text{RePair}(\text{PD}_n).$$