

Counting simple cycles

in the de Bruijn graph

Léo Ackermann¹, Pierre Peterlongo¹

word of length k

Definition (minimizer). A minimizer scheme selects within each sliding window of w k-mers over a string S, the smallest of those k-mers (breaking ties on leftmost), with respect to some order O.

Eg.

S = ATTCCTAGA

word of length k

Definition (minimizer). A minimizer scheme selects within each sliding window of w k-mers over a string S, the smallest of those k-mers (breaking ties on leftmost), with respect to some order O.

Eg.

$$S = ATTCCTAGA$$

4-mers of S

word of length k

Definition (minimizer). A minimizer scheme selects within each sliding window of w k-mers over a string S, the smallest of those k-mers (breaking ties on leftmost), with respect to some order O.

Eg.

$$S = ATTCCTAGA$$

$$ATTC$$

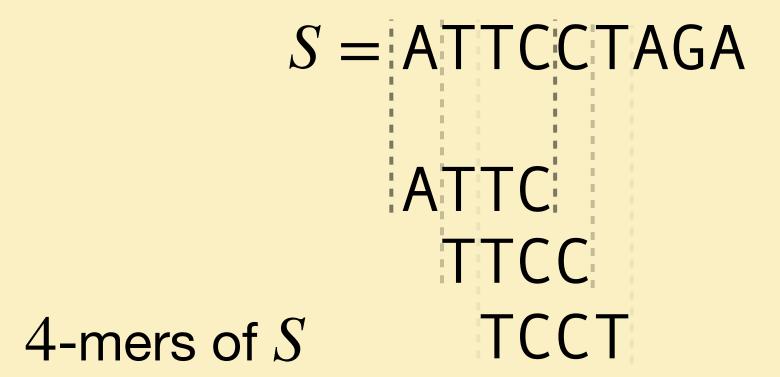
$$TTCC$$

4-mers of S

word of length k

Definition (minimizer). A minimizer scheme selects within each sliding window of w k-mers over a string S, the smallest of those k-mers (breaking ties on leftmost), with respect to some order O.

Eg.



word of length k

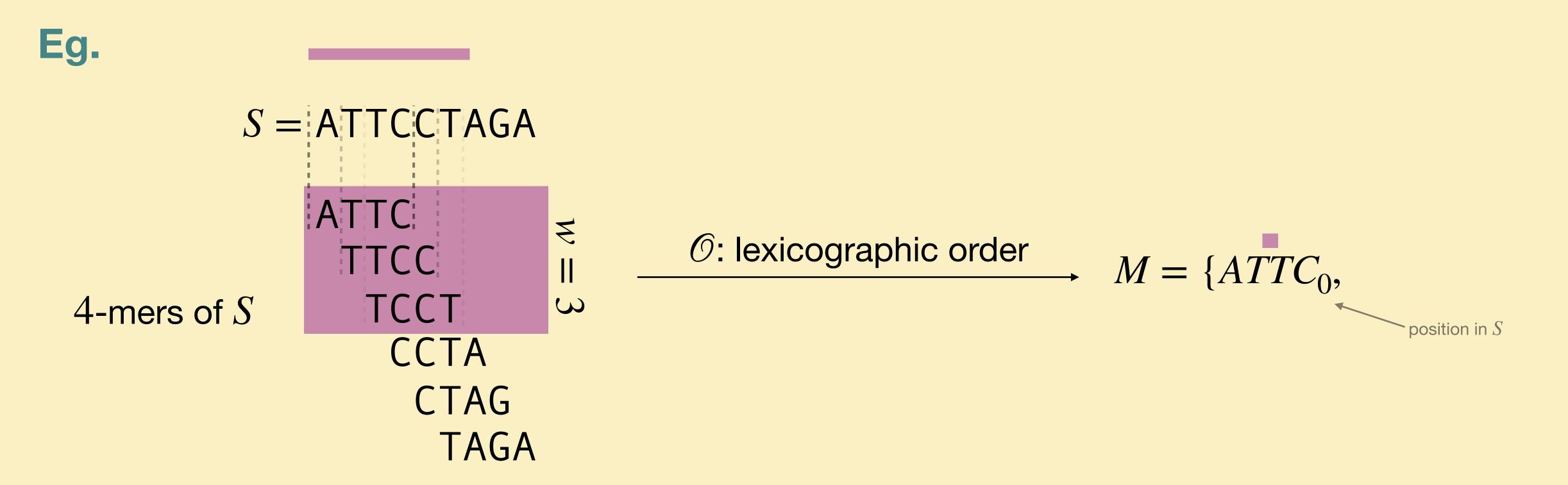
Definition (minimizer). A minimizer scheme selects within each sliding window of w k-mers over a string S, the smallest of those k-mers (breaking ties on leftmost), with respect to some order O.

Eg.

```
S = \mathsf{ATTCCTAGA}
\mathsf{ATTC}
\mathsf{TTCC}
\mathsf{4-mers} \ \mathsf{of} \ S
\mathsf{TCCT}
\mathsf{CCTA}
\mathsf{CTAG}
\mathsf{TAGA}
```

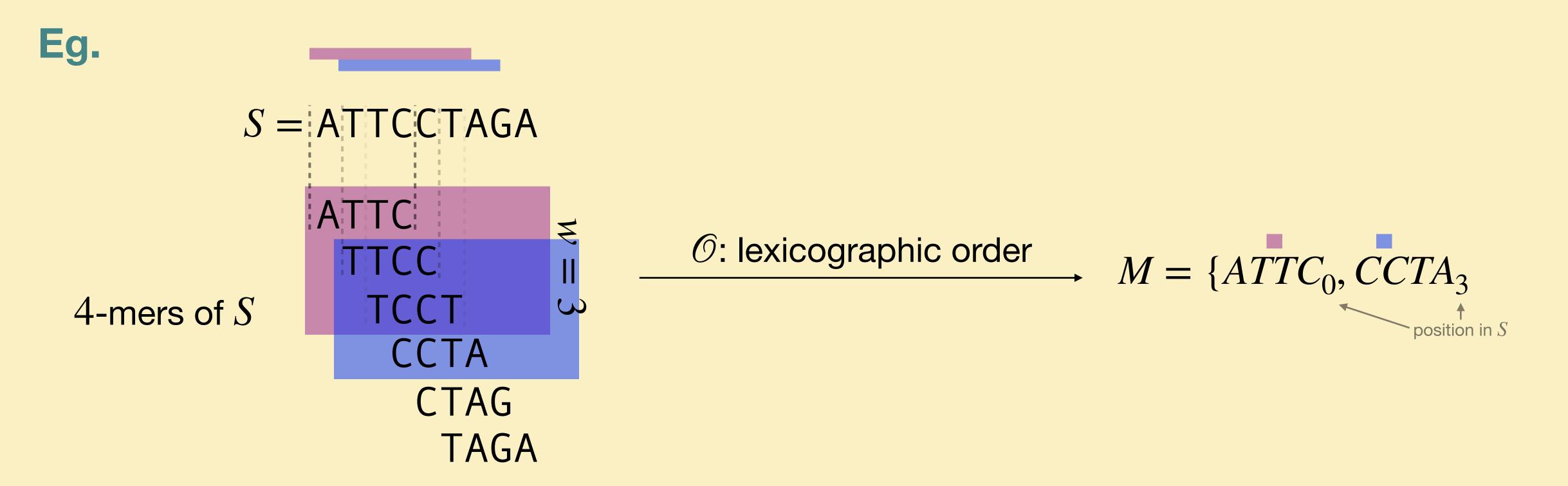
word of length k

Definition (minimizer). A minimizer scheme selects within each sliding window of w k-mers over a string S, the smallest of those k-mers (breaking ties on leftmost), with respect to some order O.



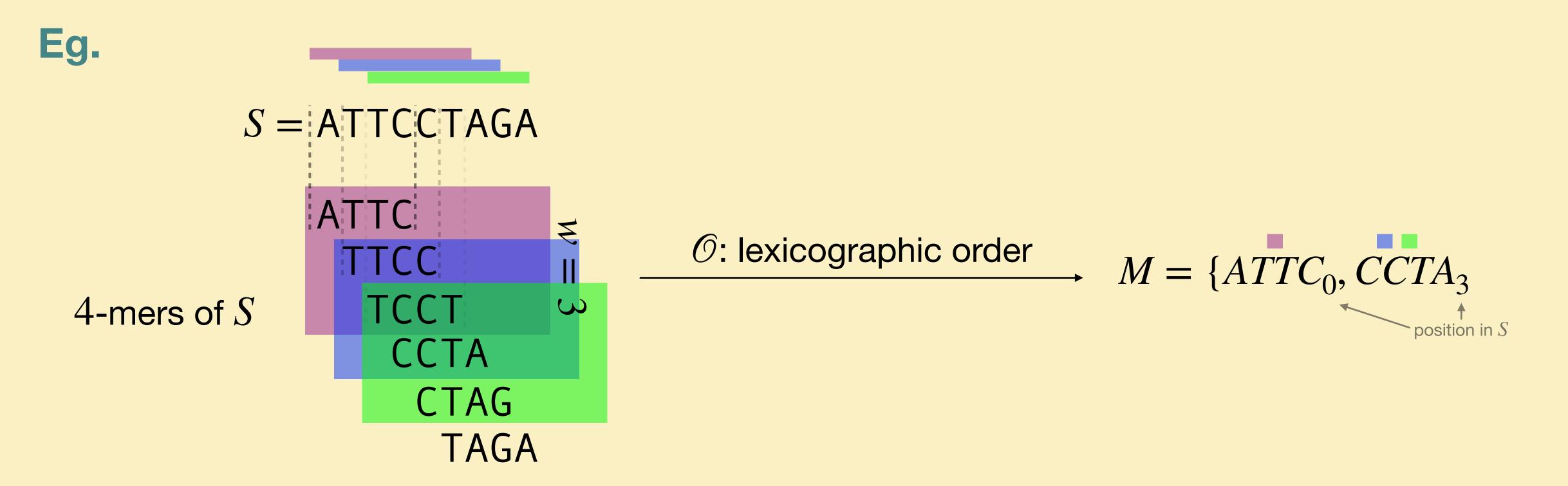
word of length k

Definition (minimizer). A minimizer scheme selects within each sliding window of w k-mers over a string S, the smallest of those k-mers (breaking ties on leftmost), with respect to some order O.



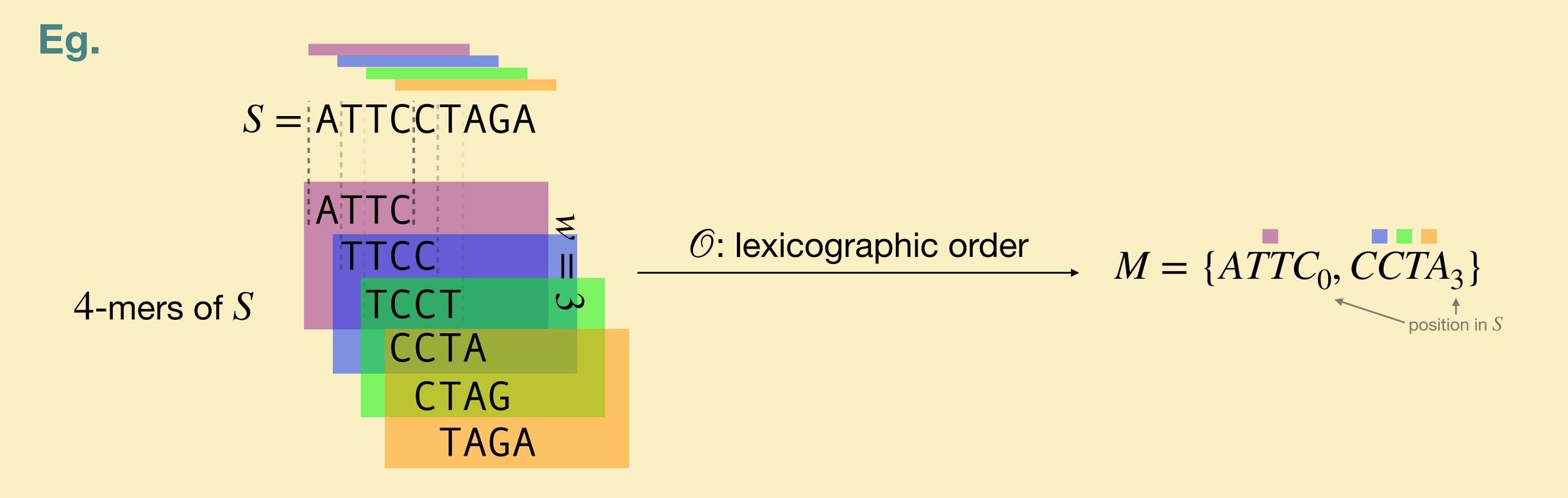
word of length k

Definition (minimizer). A minimizer scheme selects within each sliding window of w k-mers over a string S, the smallest of those k-mers (breaking ties on leftmost), with respect to some order G.



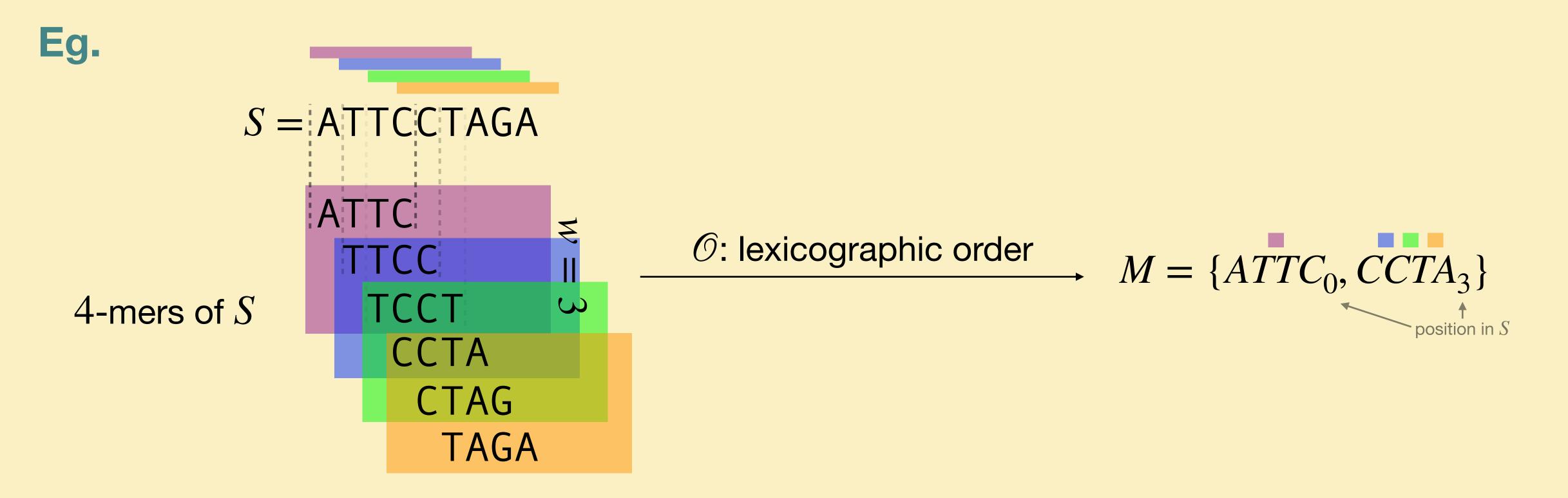
word of length k

Definition (minimizer). A minimizer scheme selects within each sliding window of w k-mers over a string S, the smallest of those k-mers (breaking ties on leftmost), with respect to some order O.



word of length k

Definition (minimizer). A minimizer scheme selects within each sliding window of w k-mers over a string S, the smallest of those k-mers (breaking ties on leftmost), with respect to some order O.

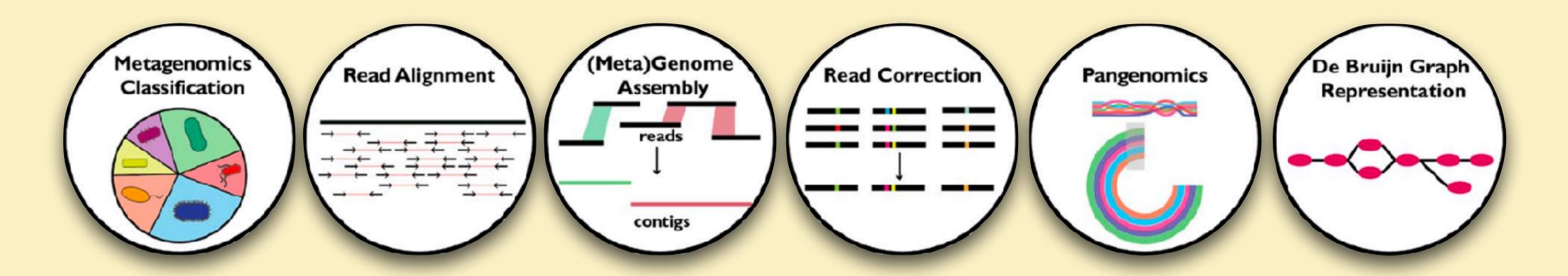


Typically, it holds that
$$|\mathcal{G}| \gg |\operatorname{sp}^k(\mathcal{G})| \gg |\operatorname{mnz}^{k,w,\mathcal{O}}(\mathcal{G})|$$
.

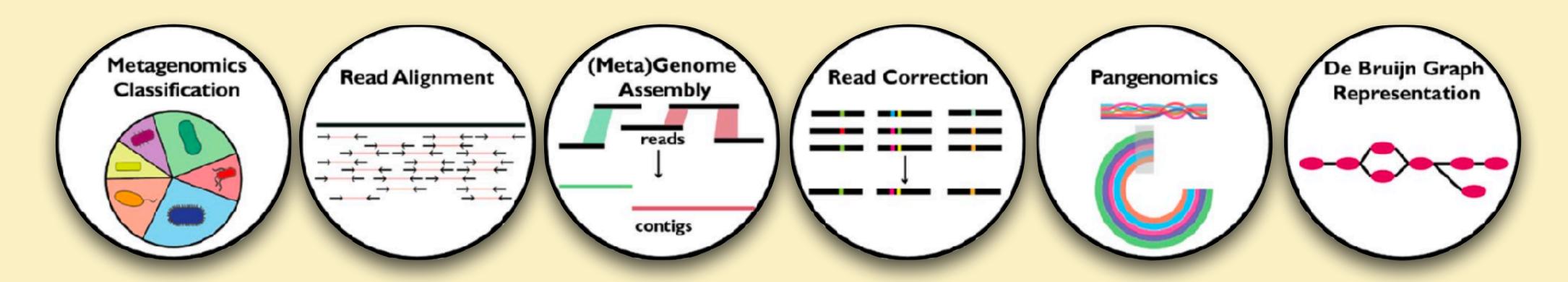
set of genomes set of k -mers set of minimizers

⇒ good for big data era

Applications of minimizers



Applications of minimizers

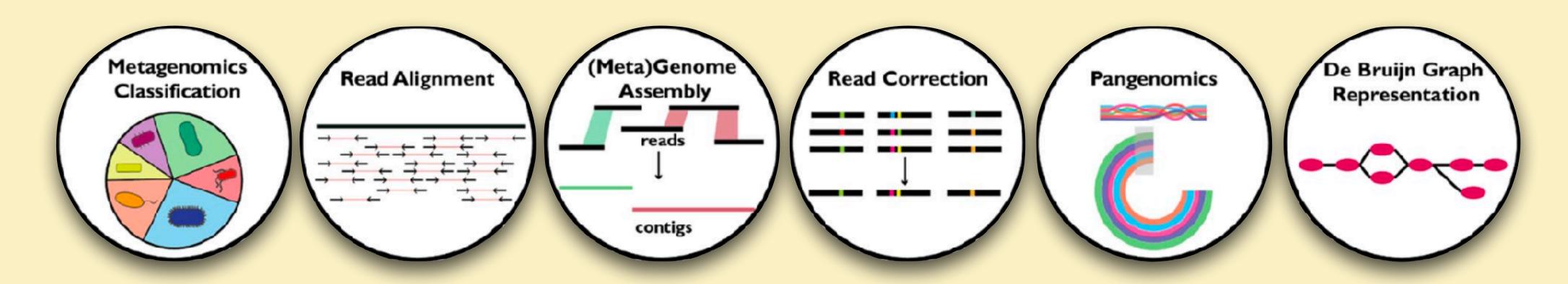


The sparser the better. The (specific) density of a minimizer scheme on a string S is

$$D = \frac{\left| \operatorname{mnz}^{k,w,\mathcal{O}}(S) \right|}{\left| S \right| - w - k + 2}.$$

It ranges somewhere within 1/w and 1.

Applications of minimizers



The sparser the better. The (specific) density of a minimizer scheme on a string S is

$$D = \frac{\left| \operatorname{mnz}^{k,w,\mathscr{O}}(S) \right|}{\left| S \right| - w - k + 2}.$$

It ranges somewhere within 1/w and 1.—every k-mer is selected

a k-mer is selected as soon as it enters the windows and as long as it lies within it

Extremal examples (on lexicographic order, k = 4, w = 3).

$$-S = CCACCACCACCC$$
 (D = $1/w$)

Definition (greedy minimizer). The tie breaking rule favors the last selected minimizer, and takes the rightmost otherwise. This greatly reduces the density.

Eg.

Minimizer

Greedy minimizer

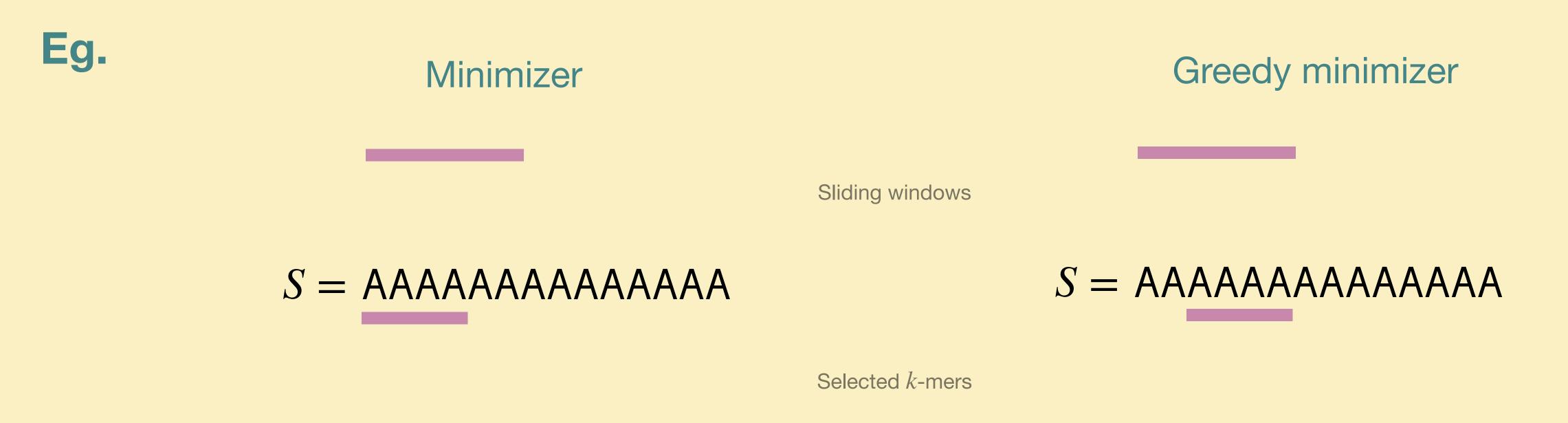
Sliding windows

S = AAAAAAAAAAAAA

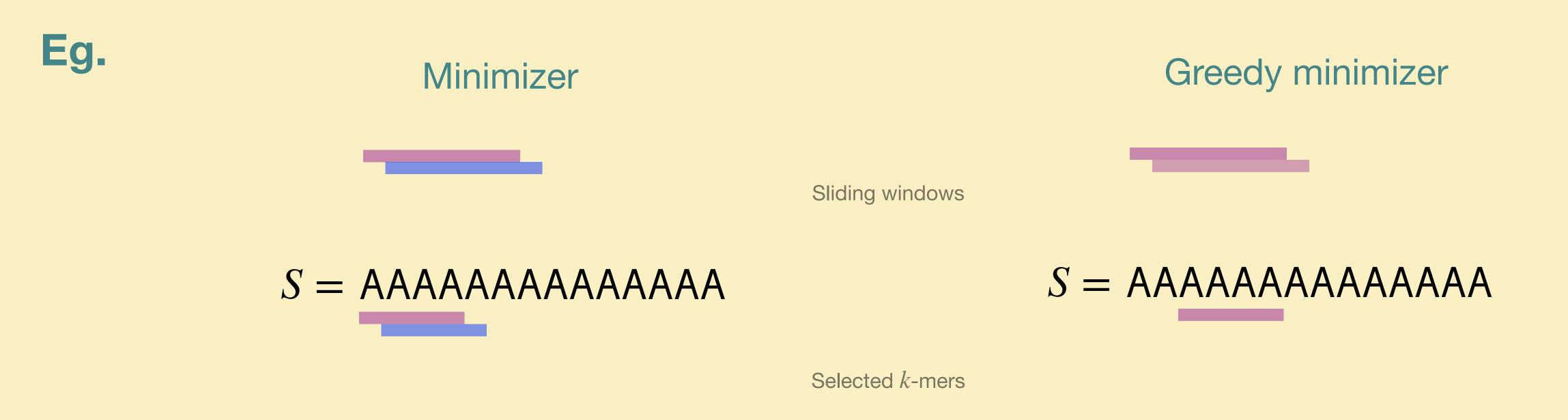
S = AAAAAAAAAAAAA

Selected *k*-mers

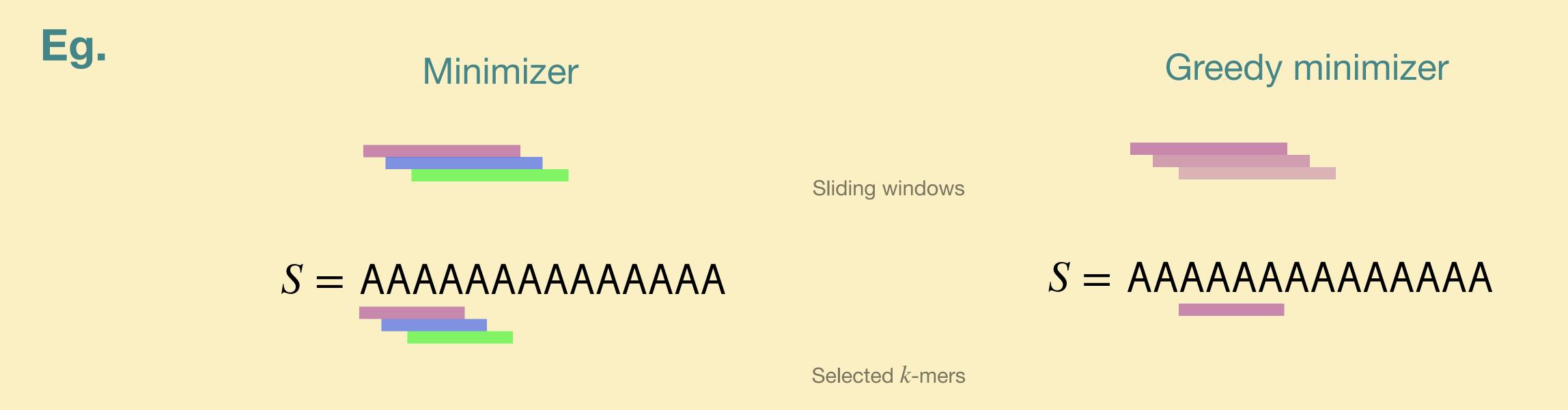
Definition (greedy minimizer). The tie breaking rule favors the last selected minimizer, and takes the rightmost otherwise. This greatly reduces the density.



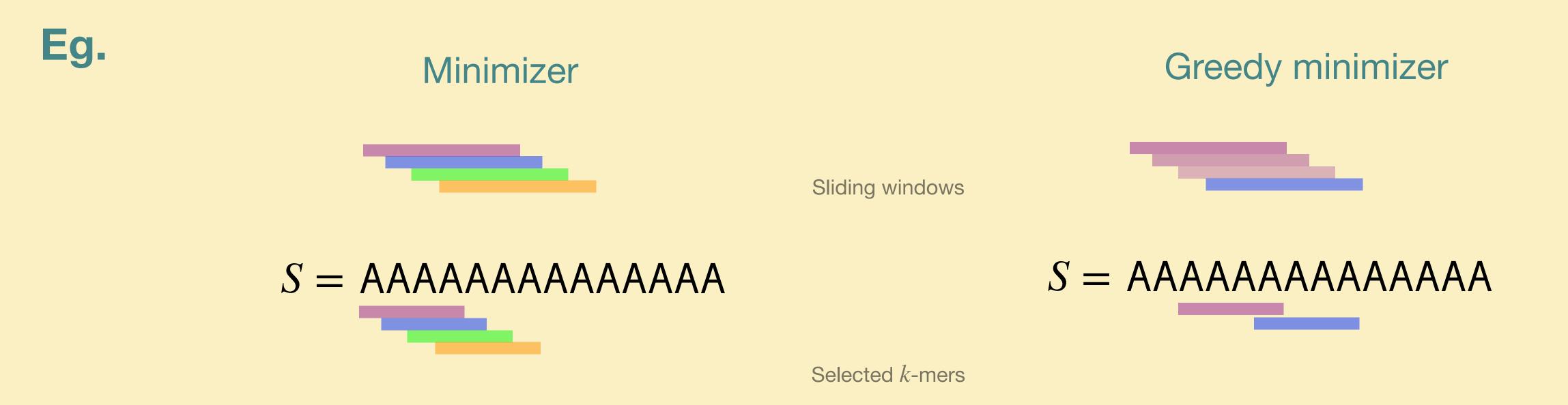
Definition (greedy minimizer). The tie breaking rule favors the last selected minimizer, and takes the rightmost otherwise. This greatly reduces the density.



Definition (greedy minimizer). The tie breaking rule favors the last selected minimizer, and takes the rightmost otherwise. This greatly reduces the density.



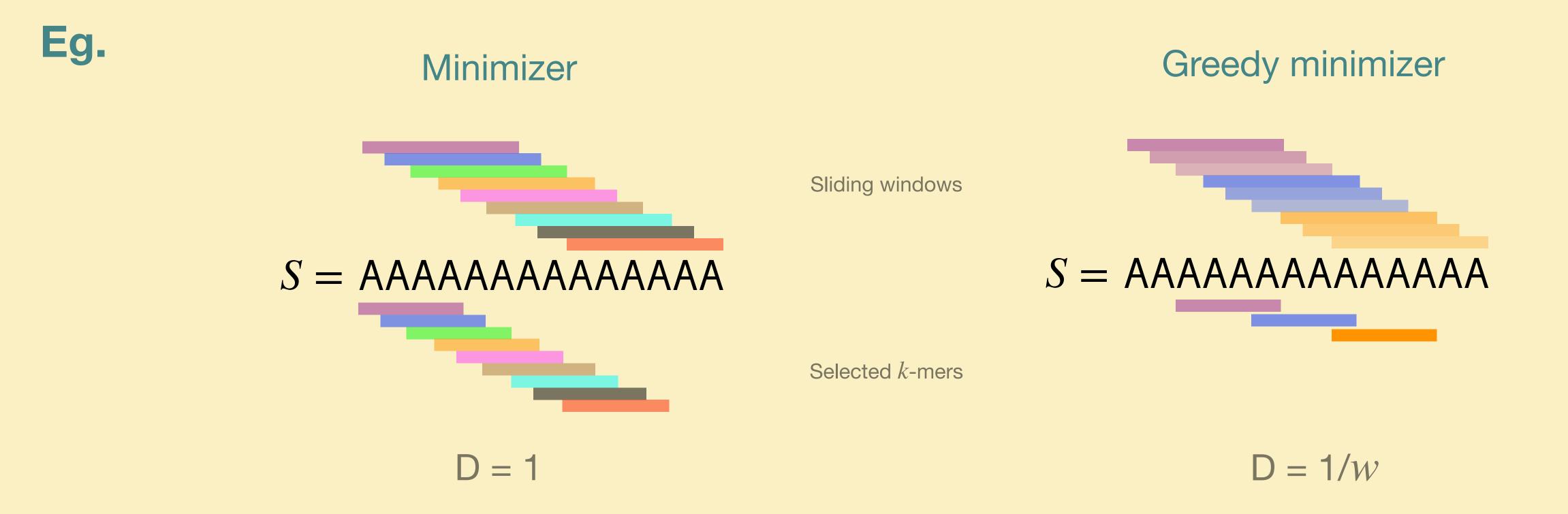
Definition (greedy minimizer). The tie breaking rule favors the last selected minimizer, and takes the rightmost otherwise. This greatly reduces the density.



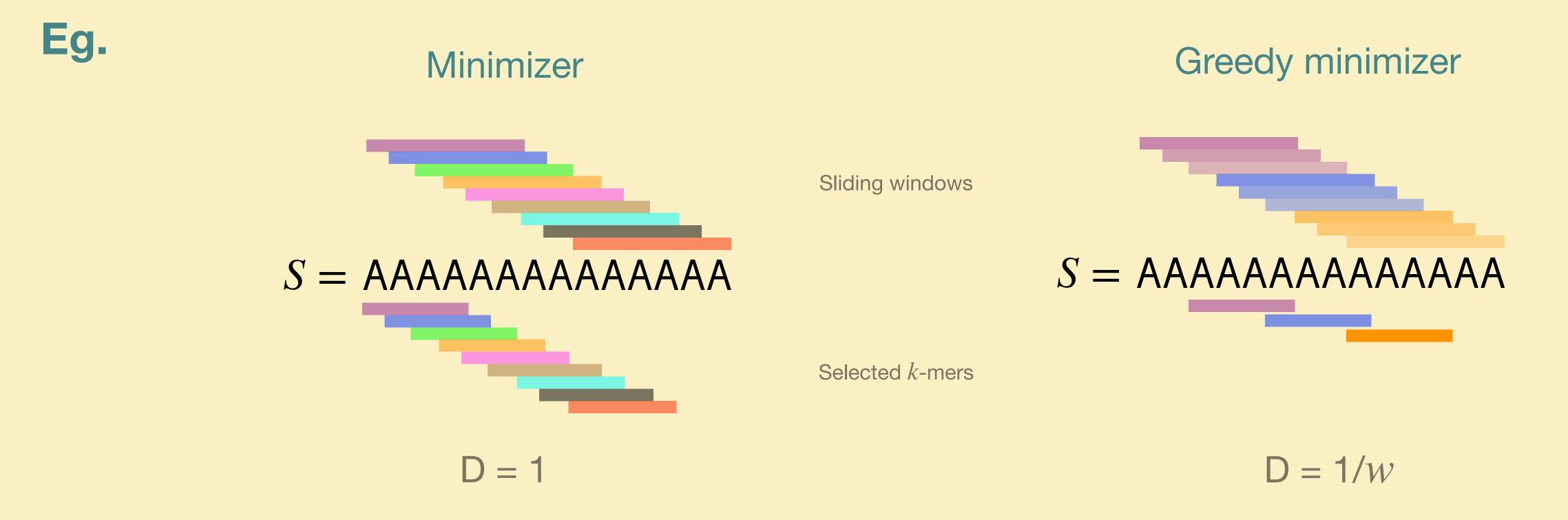
Definition (greedy minimizer). The tie breaking rule favors the last selected minimizer, and takes the rightmost otherwise. This greatly reduces the density.



Definition (greedy minimizer). The tie breaking rule favors the last selected minimizer, and takes the rightmost otherwise. This greatly reduces the density.



Definition (greedy minimizer). The tie breaking rule favors the last selected minimizer, and takes the rightmost otherwise. This greatly reduces the density.



Research question. Given w and k, come with a couple (\mathcal{O}, S) such that the specific density of the greedy minimizer (w, k, \mathcal{O}) on S is maximal.

Observation. There are two mechanisms that cause a greedy minimizer to change selected k-mer between consecutive windows:

Observation. There are two mechanisms that cause a greedy minimizer to change selected k-mer between consecutive windows:

1. The previous k-mer would have been selected (strictly smaller), but is no longer within the window,

Observation. There are two mechanisms that cause a greedy minimizer to change selected k-mer between consecutive windows:

- 1. The previous k-mer would have been selected (strictly smaller), but is no longer within the window,
- 2. The new k-mer, appearing at the right of the window, is the (strictly) smallest one.

Observation. There are two mechanisms that cause a greedy minimizer to change selected k-mer between consecutive windows:

- 1. The previous k-mer would have been selected (strictly smaller), but is no longer within the window,
- 2. The new k-mer, appearing at the right of the window, is the (strictly) smallest one.

A maximal density string.



Observation. There are two mechanisms that cause a greedy minimizer to change selected k-mer between consecutive windows:

- 1. The previous k-mer would have been selected (strictly smaller), but is no longer within the window,
- 2. The new k-mer, appearing at the right of the window, is the (strictly) smallest one.

A maximal density string.

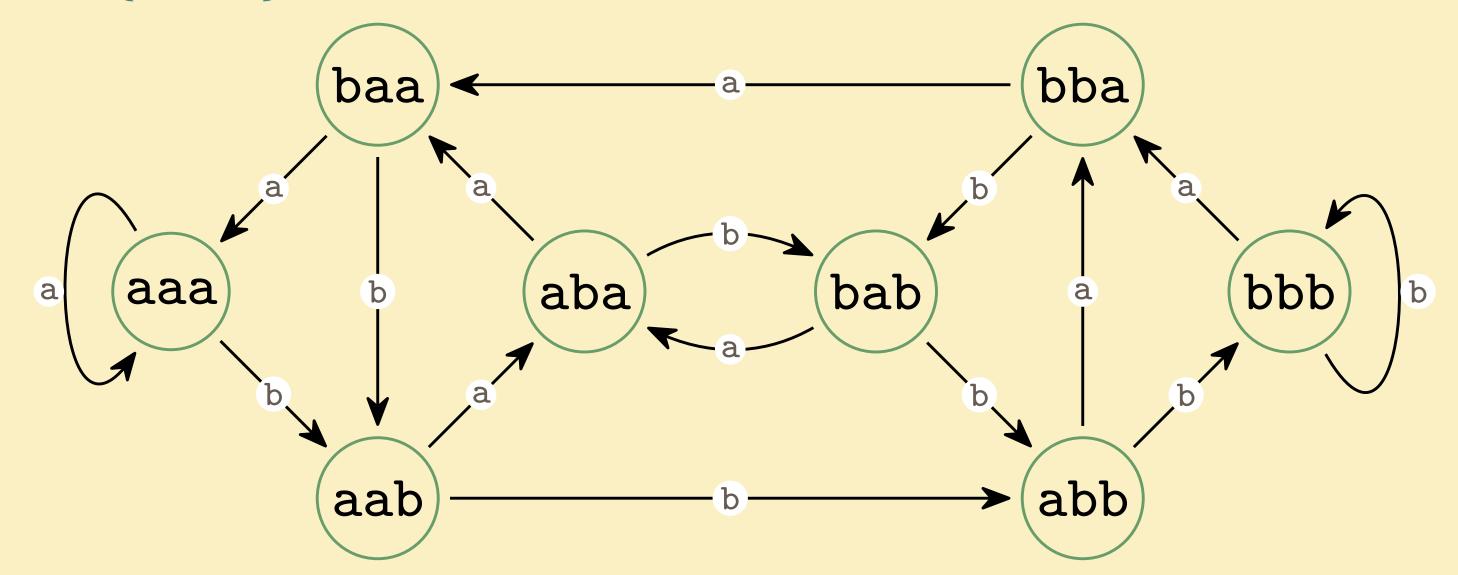


Research question (rough reformulation). Come with an ordering on the vertices of the de Bruijn graph of order k such that there exist an increasing-then-decreasing path as long as possible in this dBG.

Definition (de Bruijn graph). The (complete) dBG of order n over an alphabet Σ is the graph whose set of vertices is Σ^n , and edges are given by: $u \to v$ iff the (n-1)-long suffix of u is a prefix of v.

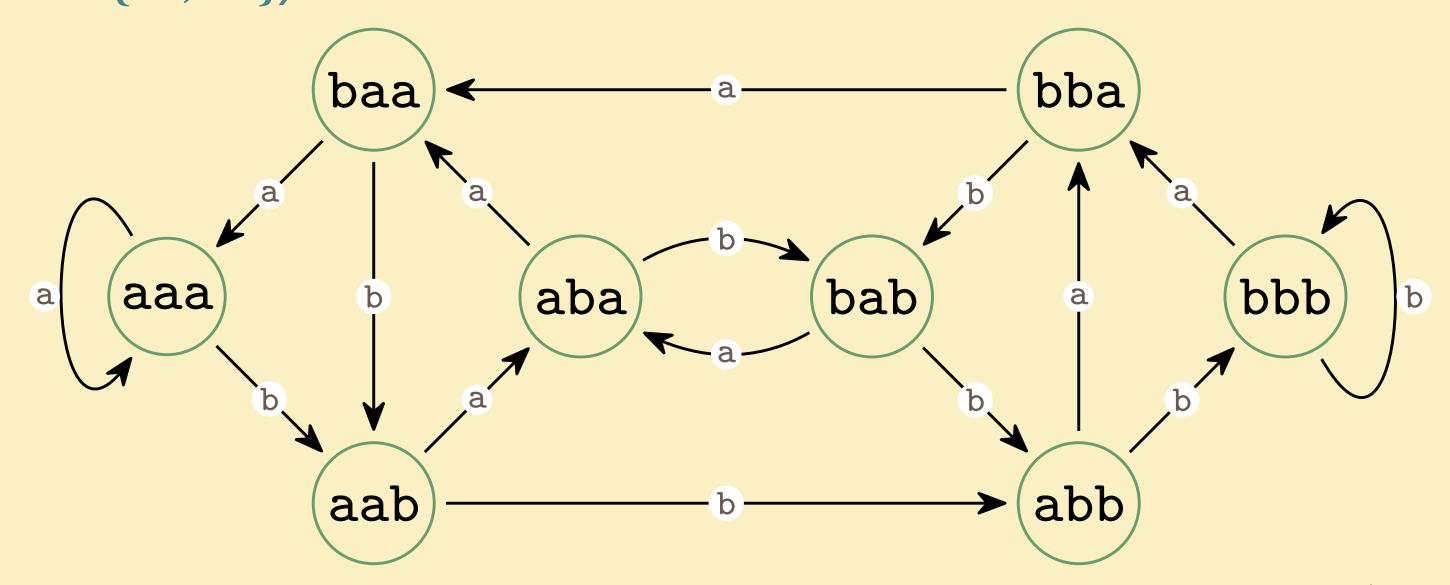
Definition (de Bruijn graph). The (complete) dBG of order n over an alphabet Σ is the graph whose set of vertices is Σ^n , and edges are given by: $u \to v$ iff the (n-1)-long suffix of u is a prefix of v.

Eg.
$$dBG(n = 3, \Sigma = \{a, b\})$$
.



Definition (de Bruijn graph). The (complete) dBG of order n over an alphabet Σ is the graph whose set of vertices is Σ^n , and edges are given by: $u \to v$ iff the (n-1)-long suffix of u is a prefix of v.

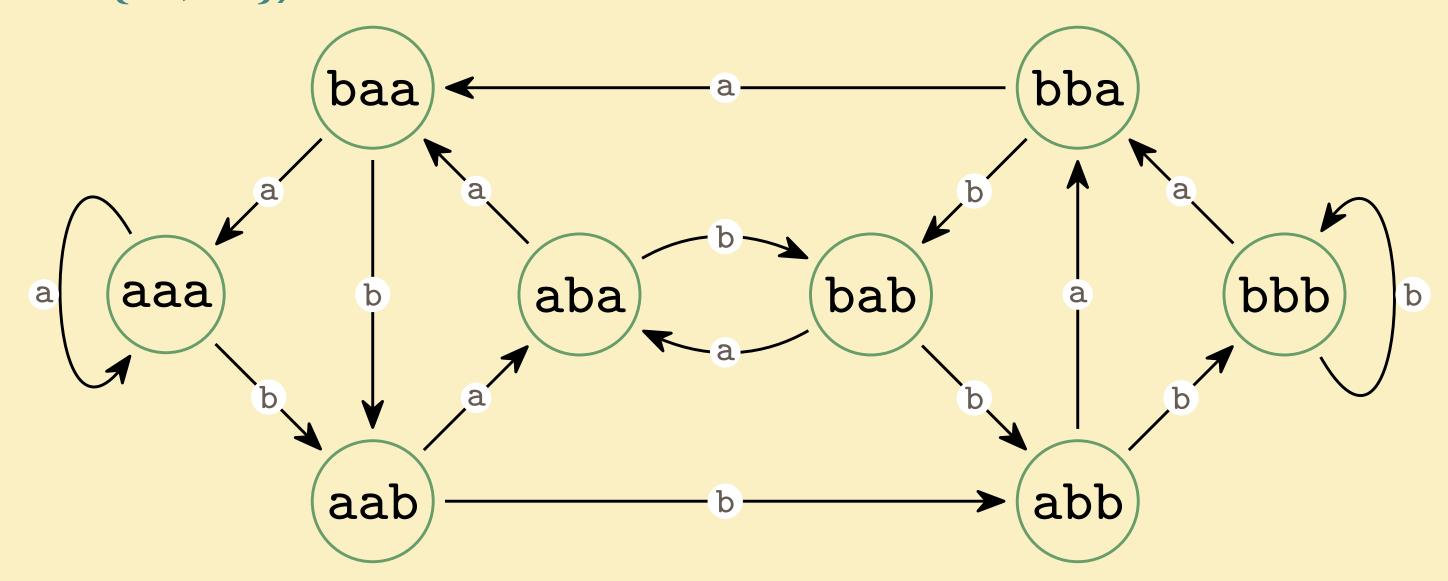
Eg.
$$dBG(n = 3, \Sigma = \{a, b\})$$
.



Combinatorial explosion. These graphs have $|\Sigma|^n$ vertices, and $|\Sigma|^{n+1}$ edges.

Definition (de Bruijn graph). The (complete) dBG of order n over an alphabet Σ is the graph whose set of vertices is Σ^n , and edges are given by: $u \to v$ iff the (n-1)-long suffix of u is a prefix of v.

Eg.
$$dBG(n = 3, \Sigma = \{a, b\})$$
.

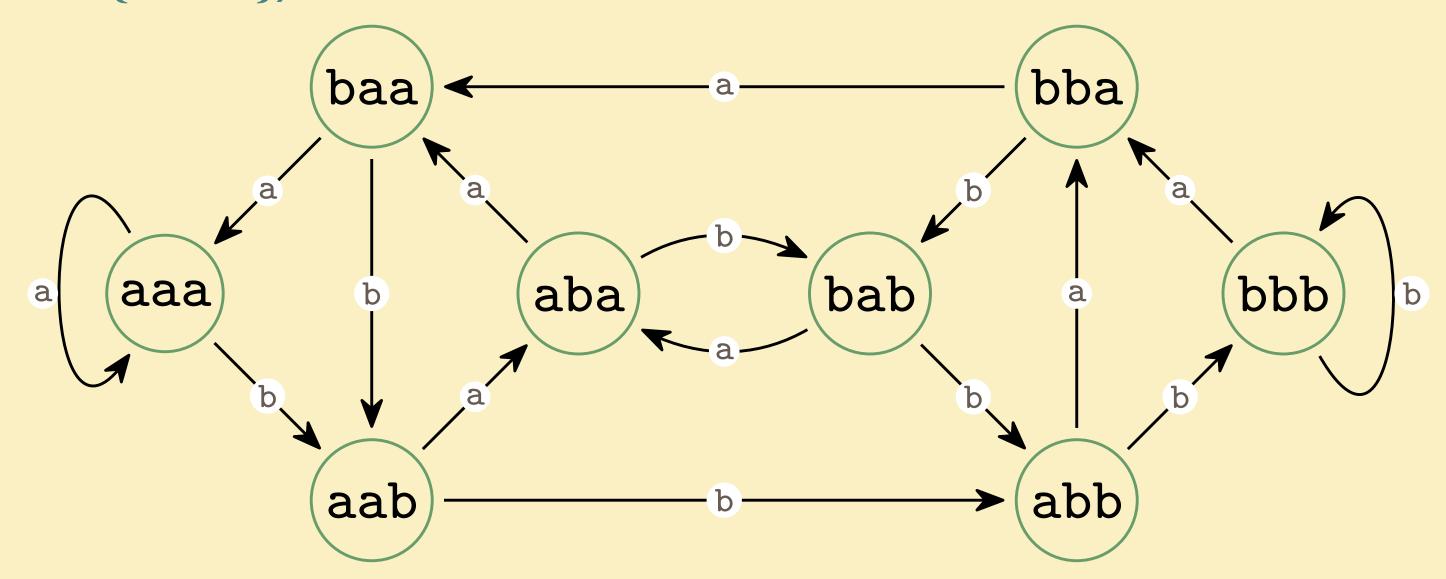


Combinatorial explosion. These graphs have $|\Sigma|^n$ vertices, and $|\Sigma|^{n+1}$ edges.

 \Rightarrow cannot be explored computationally (for typical bioinformatics values: $4^{31} \approx 4 \cdot 10^{18}$ vertices, twice as much edges)

Definition (de Bruijn graph). The (complete) dBG of order n over an alphabet Σ is the graph whose set of vertices is Σ^n , and edges are given by: $u \to v$ iff the (n-1)-long suffix of u is a prefix of v.

Eg.
$$dBG(n = 3, \Sigma = \{a, b\})$$
.



Combinatorial explosion. These graphs have $|\Sigma|^n$ vertices, and $|\Sigma|^{n+1}$ edges.

 \Rightarrow cannot be explored computationally (for typical bioinformatics values: $4^{31} \approx 4 \cdot 10^{18}$ vertices, twice as much edges)

nterest. Sliding windows (minimizers, compression schemes)

Counting simple cycles in the de Bruijn graph

Question. How many simple cycles of length ℓ live within the *de Bruijn* graph of order n over an alphabet made of σ letters?

Counting simple cycles in the de Bruijn graph

Question. How many simple cycles of length ℓ live within the *de Bruijn* graph of order n over an alphabet made of σ letters?

Outline.

- A bijection from simple cycles to a restriction of Lyndon words
- Reformulation of the state of knowledge
- Extending knowledge on simple cases (+ conjecture)
- A small CLI interface to count/enumerate simple cycles for the curious

Correspondence with restricted class of Lyndon words

Definition (Lyndon word). A Lyndon word is a word strictly lexicographically smaller than its rotations.

Definition (Lyndon word). A Lyndon word is a word strictly lexicographically smaller than its rotations.

Proposition. There is an explicit formula for $\lambda_{\sigma}(\ell)$, the number of Lyndon words of length ℓ over an alphabet of size σ .

Definition (Lyndon word). A Lyndon word is a word strictly lexicographically smaller than its rotations.

Proposition. There is an explicit formula for $\lambda_{\sigma}(\ell)$, the number of Lyndon words of length ℓ over an alphabet of size σ .

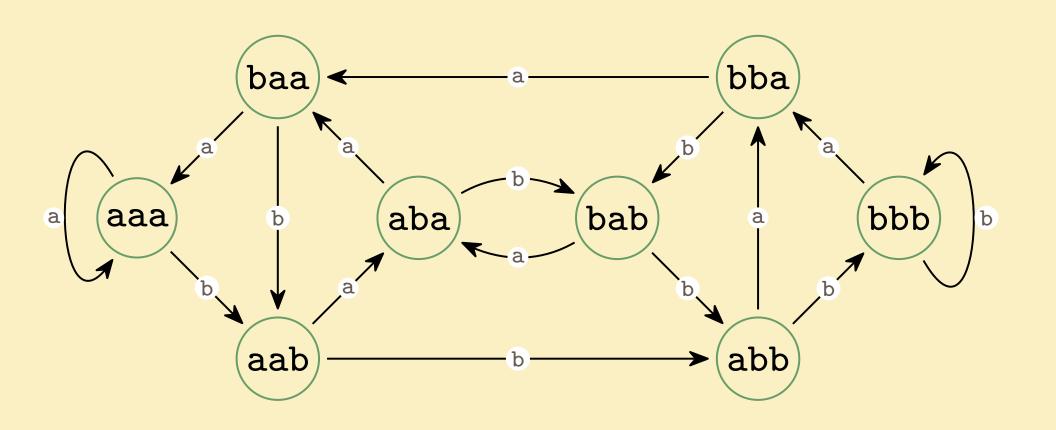
Theorem (Lempel, 1970). There is a one-to-one correspondence between Lyndon words and dBG cycles.

Eg.

Lyndon word

Powering

Sliding window



Definition (Lyndon word). A Lyndon word is a word strictly lexicographically smaller than its rotations.

Proposition. There is an explicit formula for $\lambda_{\sigma}(\ell)$, the number of Lyndon words of length ℓ over an alphabet of size σ .

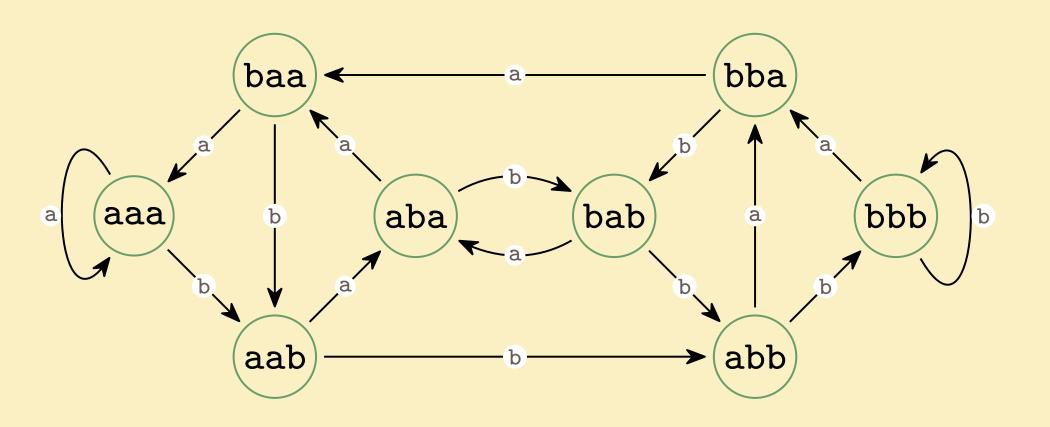
Theorem (Lempel, 1970). There is a one-to-one correspondence between Lyndon words and dBG cycles.

Eg.

Lyndon word ab

Powering

Sliding window



Definition (Lyndon word). A Lyndon word is a word strictly lexicographically smaller than its rotations.

Proposition. There is an explicit formula for $\lambda_{\sigma}(\ell)$, the number of Lyndon words of length ℓ over an alphabet of size σ .

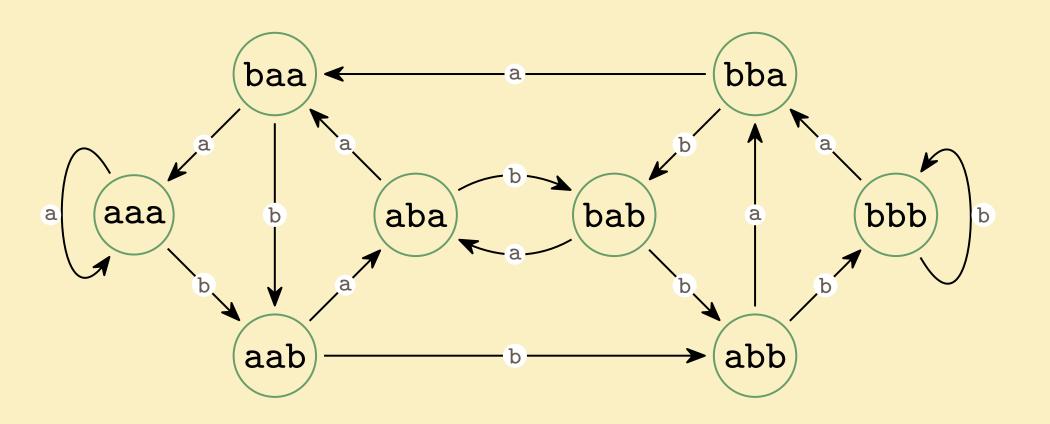
Theorem (Lempel, 1970). There is a one-to-one correspondence between Lyndon words and dBG cycles.

Eg.

Lyndon word ab

Powering abababab...

Sliding window



Definition (Lyndon word). A Lyndon word is a word strictly lexicographically smaller than its rotations.

Proposition. There is an explicit formula for $\lambda_{\sigma}(\ell)$, the number of Lyndon words of length ℓ over an alphabet of size σ .

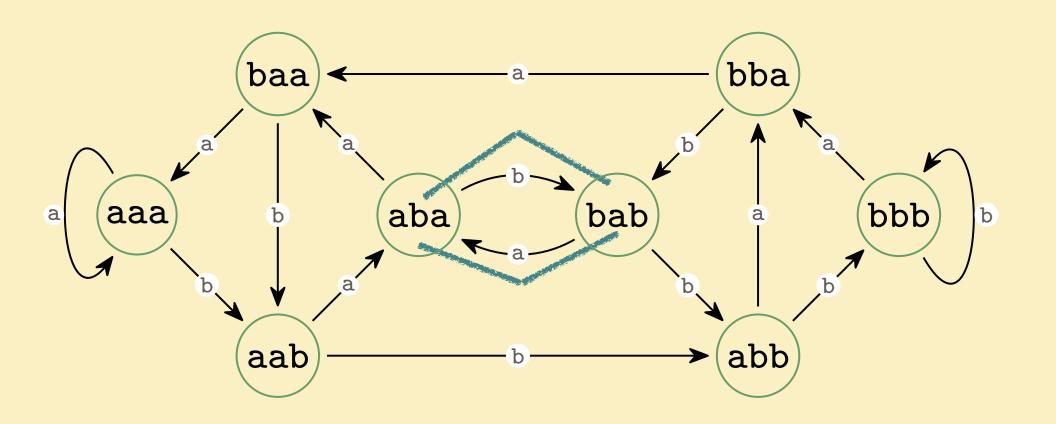
Theorem (Lempel, 1970). There is a one-to-one correspondence between Lyndon words and dBG cycles.

Eg.

Lyndon word ab

Powering abababab...

Sliding window aba-bab-aba



Definition (Lyndon word). A Lyndon word is a word strictly lexicographically smaller than its rotations.

Proposition. There is an explicit formula for $\lambda_{\sigma}(\ell)$, the number of Lyndon words of length ℓ over an alphabet of size σ .

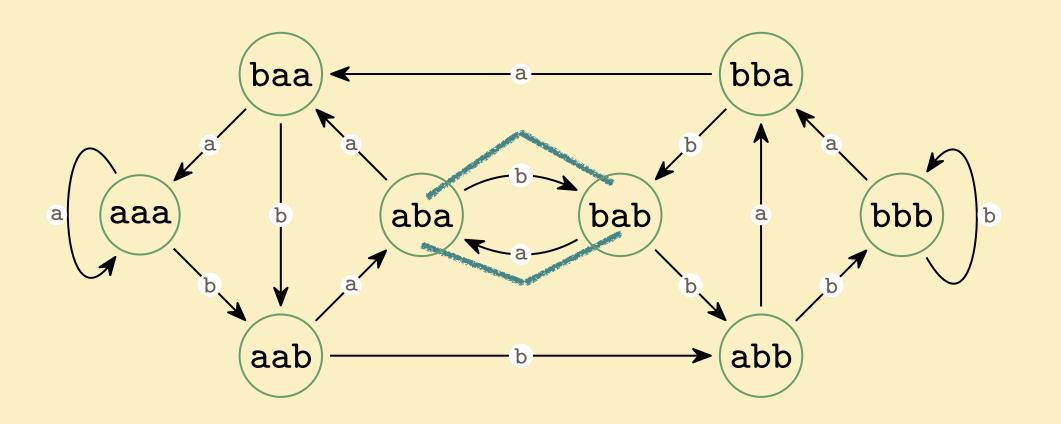
Theorem (Lempel, 1970). There is a one-to-one correspondence between Lyndon words and dBG cycles.

Eg.

Lyndon word ab aab

Powering abababab...

Sliding window aba-bab-aba



Definition (Lyndon word). A Lyndon word is a word strictly lexicographically smaller than its rotations.

Proposition. There is an explicit formula for $\lambda_{\sigma}(\ell)$, the number of Lyndon words of length ℓ over an alphabet of size σ .

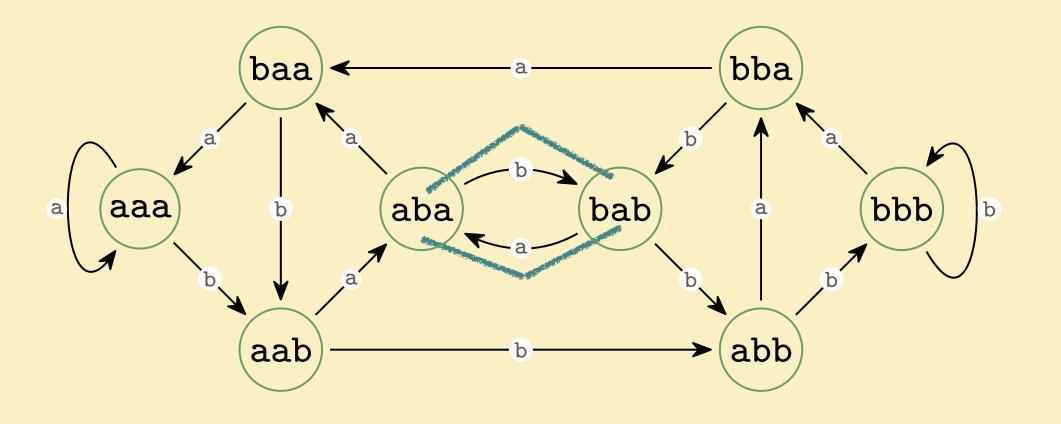
Theorem (Lempel, 1970). There is a one-to-one correspondence between Lyndon words and dBG cycles.

Eg.

Lyndon word ab aab

Powering abababab... aabaabaab...

Sliding window aba-bab-aba



Definition (Lyndon word). A Lyndon word is a word strictly lexicographically smaller than its rotations.

Proposition. There is an explicit formula for $\lambda_{\sigma}(\ell)$, the number of Lyndon words of length ℓ over an alphabet of size σ .

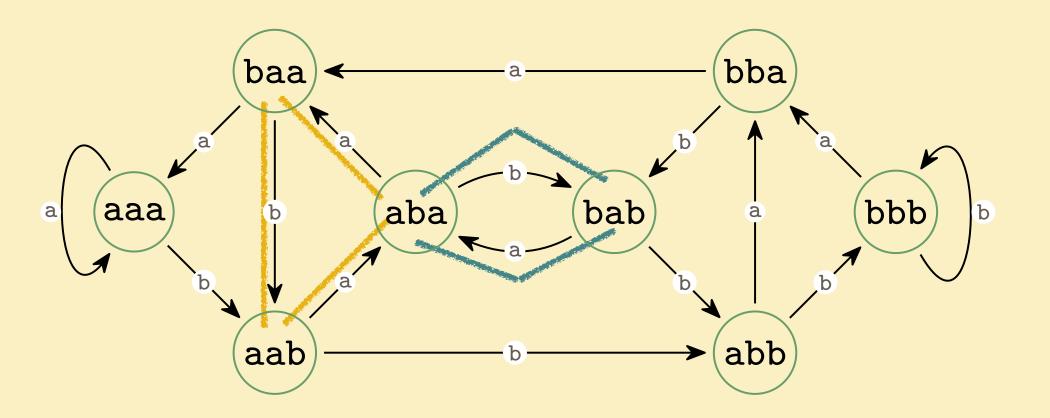
Theorem (Lempel, 1970). There is a one-to-one correspondence between Lyndon words and dBG cycles.

Eg.

Lyndon word ab aab

Powering abababab... aabaabaab...

Sliding window aba-bab-aba aab-aba-baa-aab



k-perfect Lyndon words correspond to dBG simple cycles

Definition (k-perfect Lyndon word). This is a Lyndon word whose k-mers are all distinct, when seen as a circular word.

k-perfect Lyndon words correspond to dBG simple cycles

Definition (k-perfect Lyndon word). This is a Lyndon word whose k-mers are all distinct, when seen as a circular word.

Corollary. There is a one-to-one correspondence between perfect Lyndon words and simple dBG cycles.

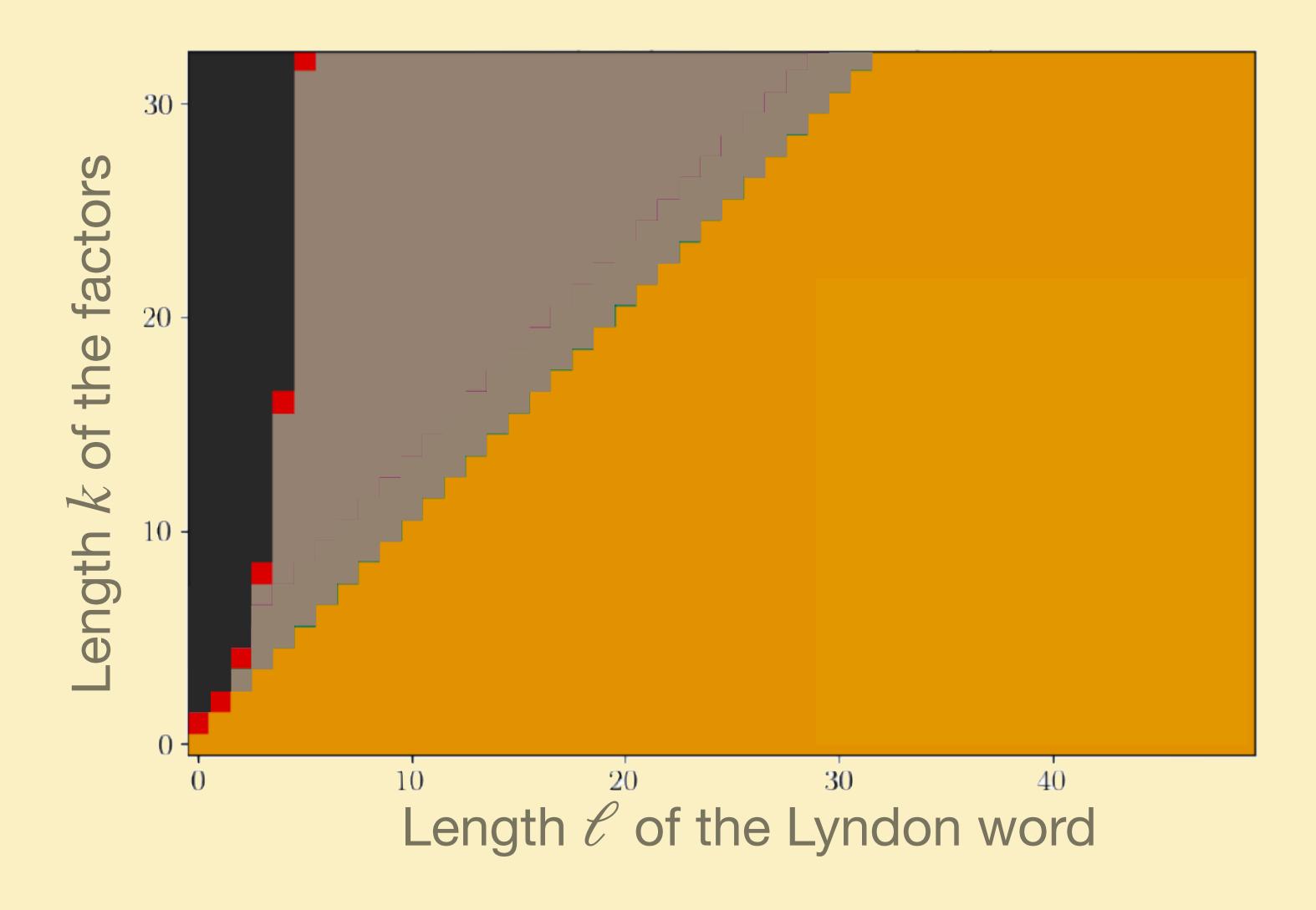
k-perfect Lyndon words correspond to dBG simple cycles

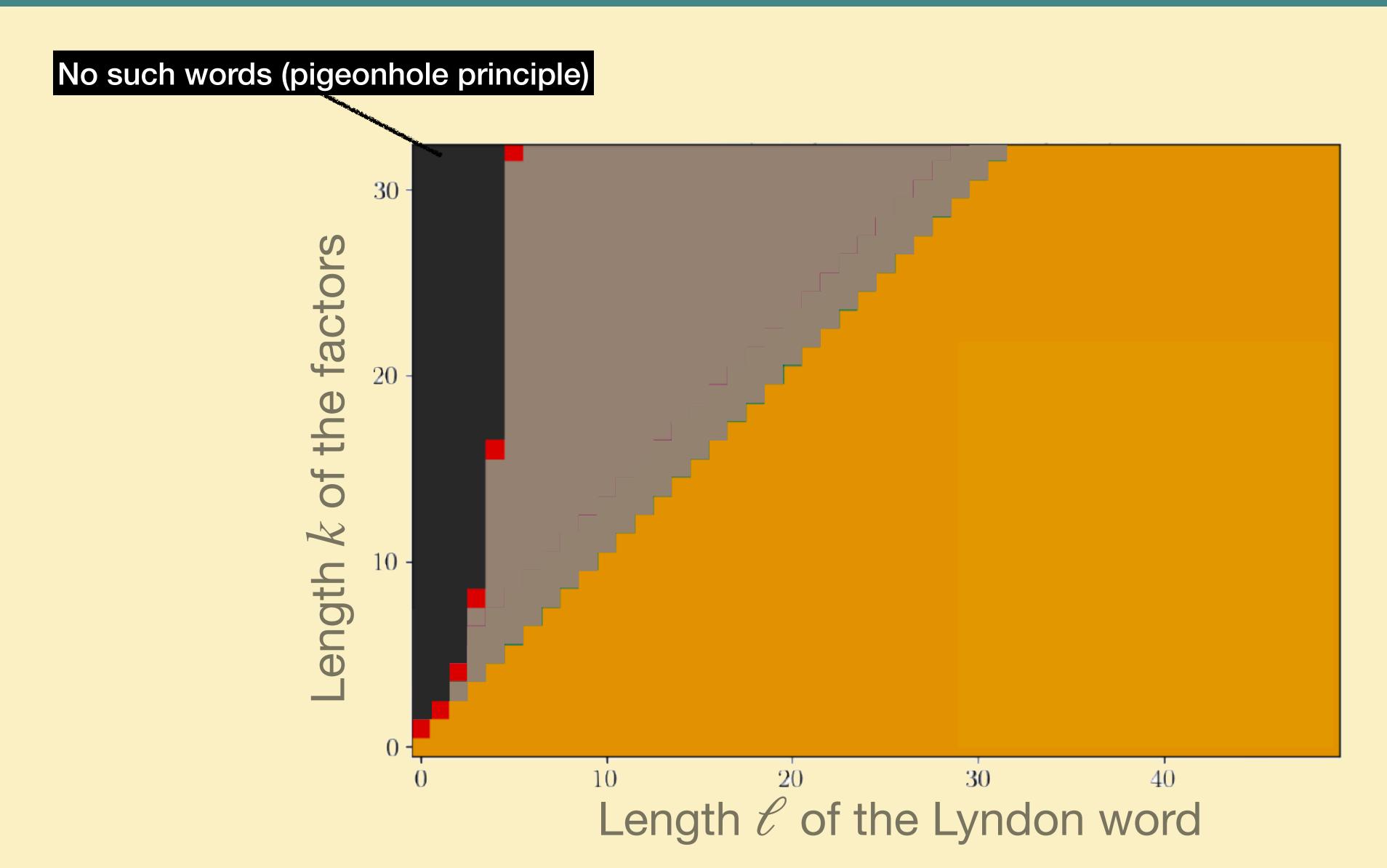
Definition (k-perfect Lyndon word). This is a Lyndon word whose k-mers are all distinct, when seen as a circular word.

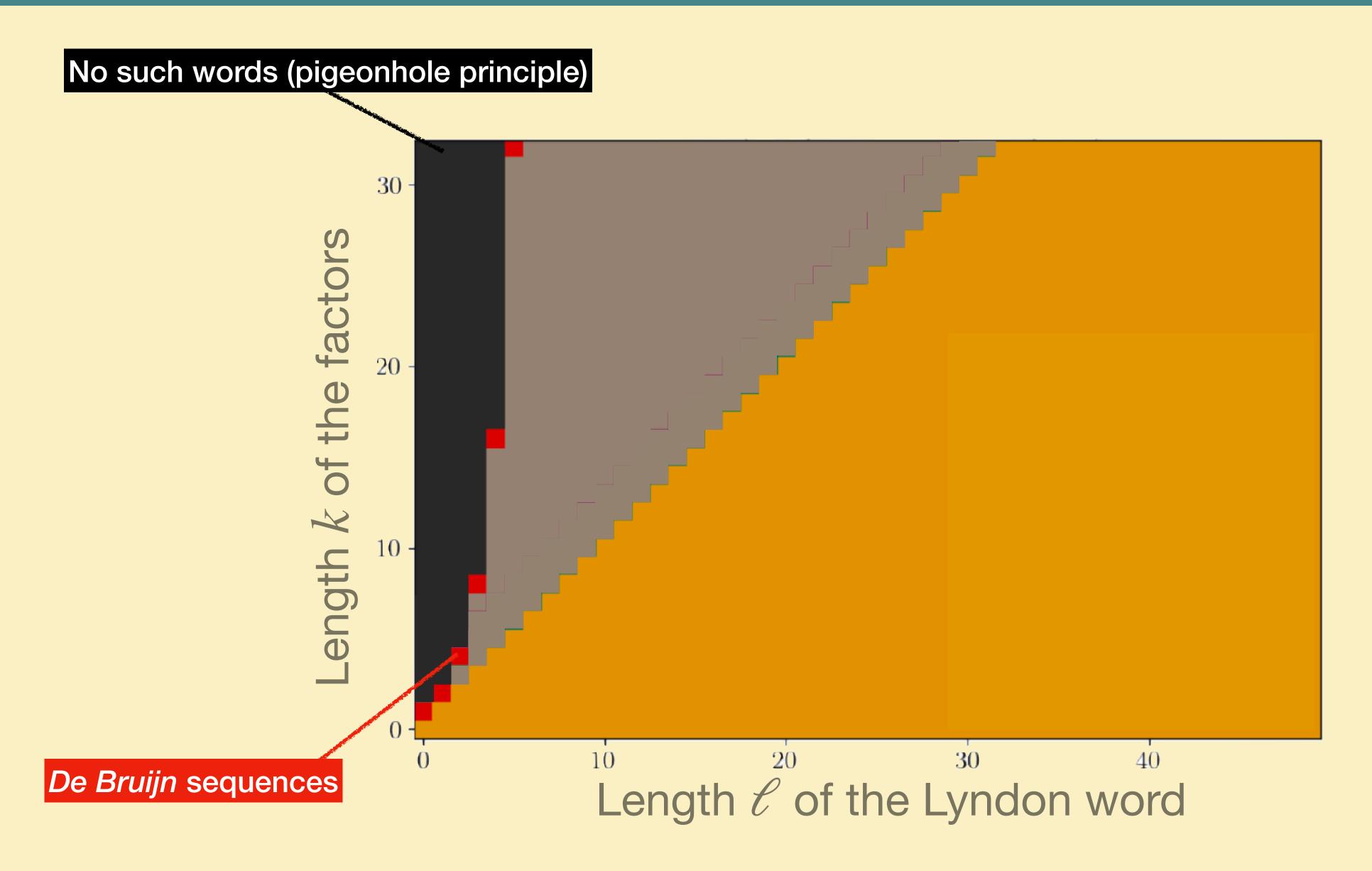
Corollary. There is a one-to-one correspondence between perfect Lyndon words and simple dBG cycles.

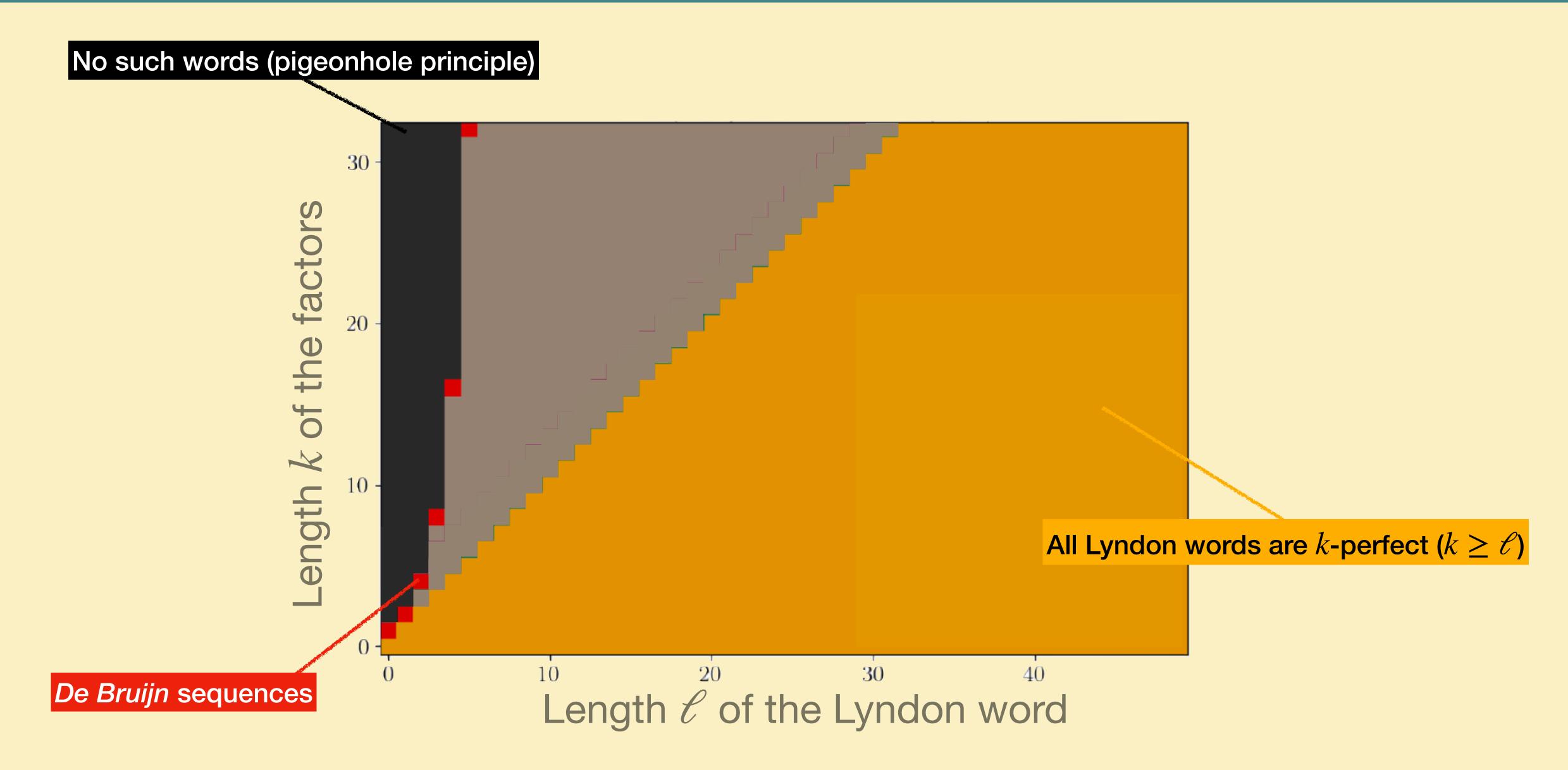
Question. How many k-perfect Lyndon words of length ℓ are there?

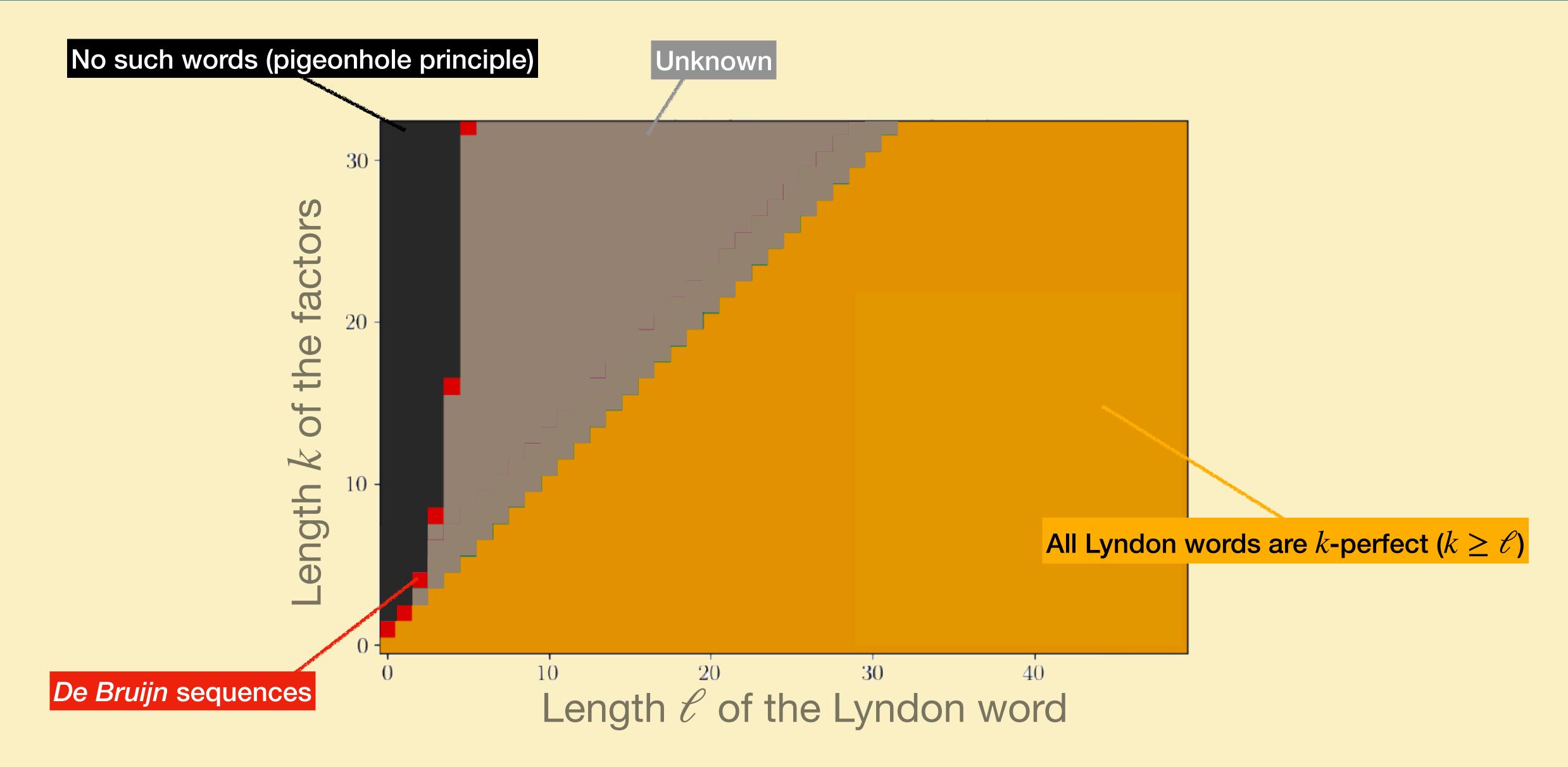
Counting simple cycles

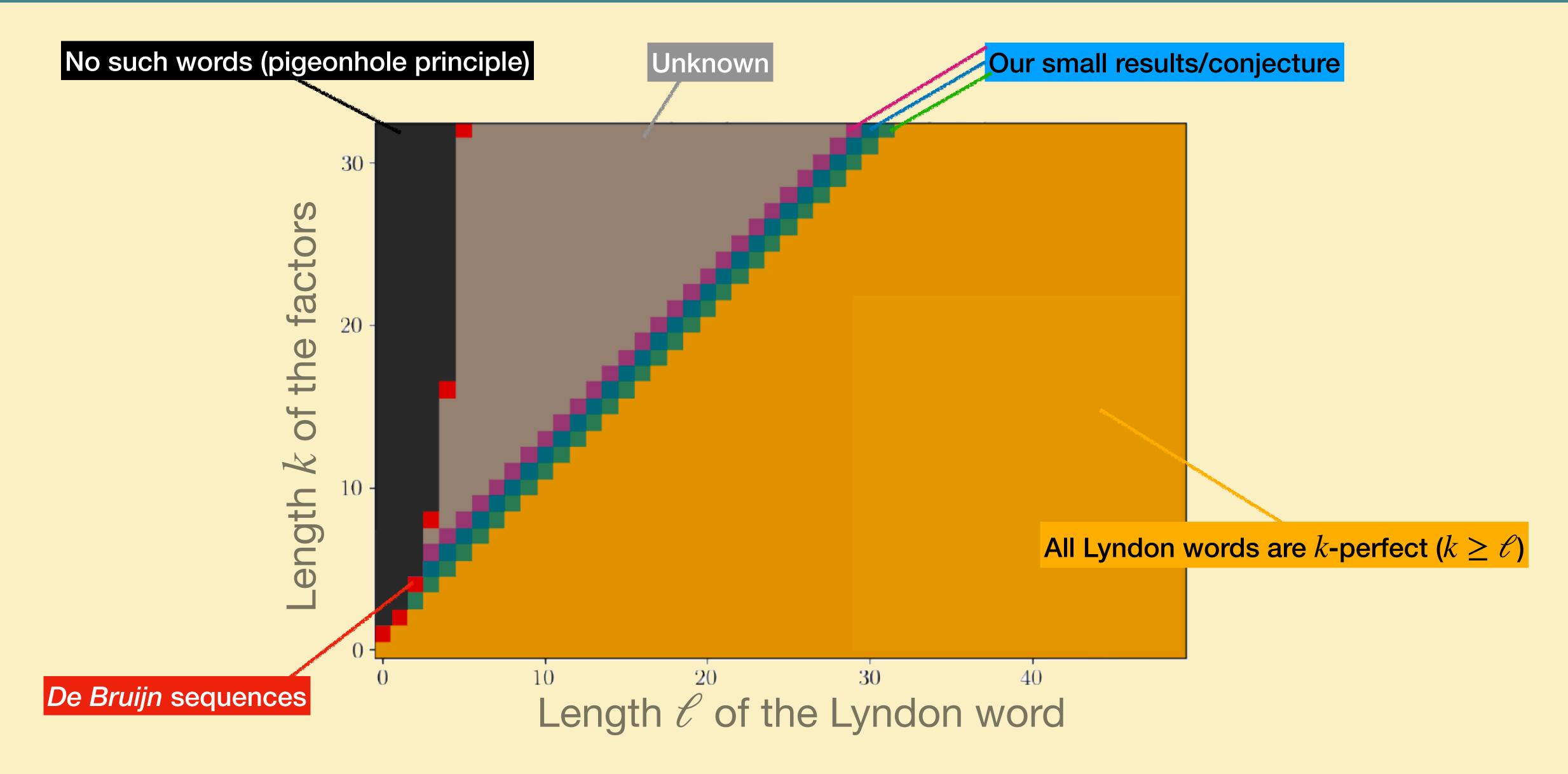












Proposition. A word w of length k+1 that admits a (circular) repeating factor of length k, necessarily has a border.

Proposition. A word w of length k+1 that admits a (circular) repeating factor of length k, necessarily has a border.

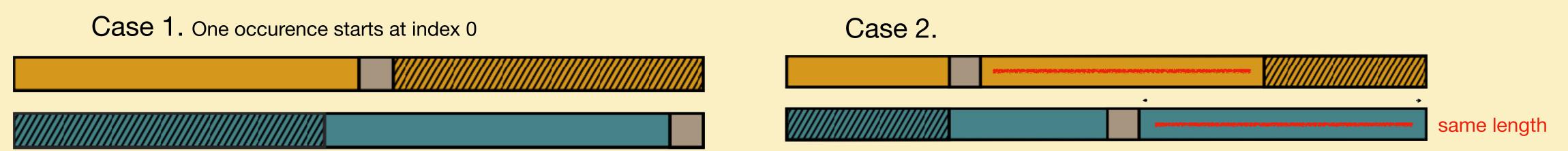
Proof.





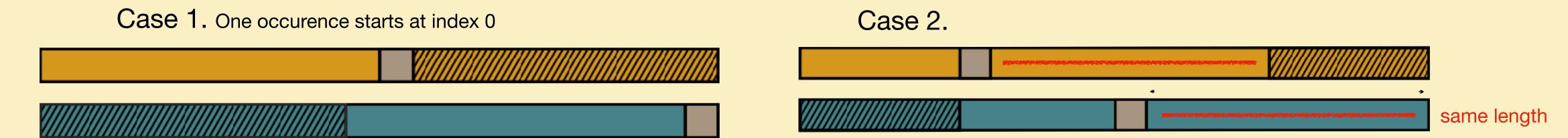
Proposition. A word w of length k+1 that admits a (circular) repeating factor of length k, necessarily has a border.

Proof.



Proposition. A word w of length k+1 that admits a (circular) repeating factor of length k, necessarily has a border.

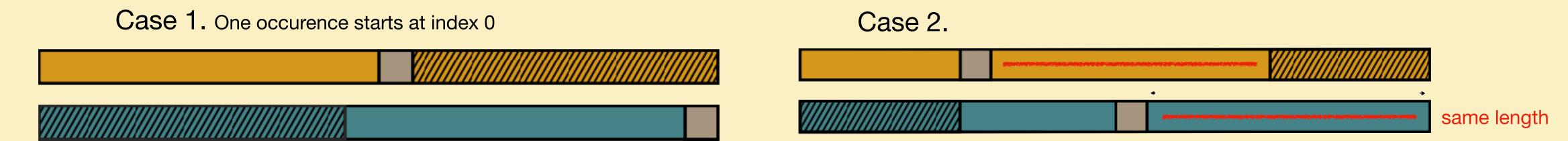
Proof.



Consequently, such a word cannot be Lyndon.

Proposition. A word w of length k+1 that admits a (circular) repeating factor of length k, necessarily has a border.

Proof.



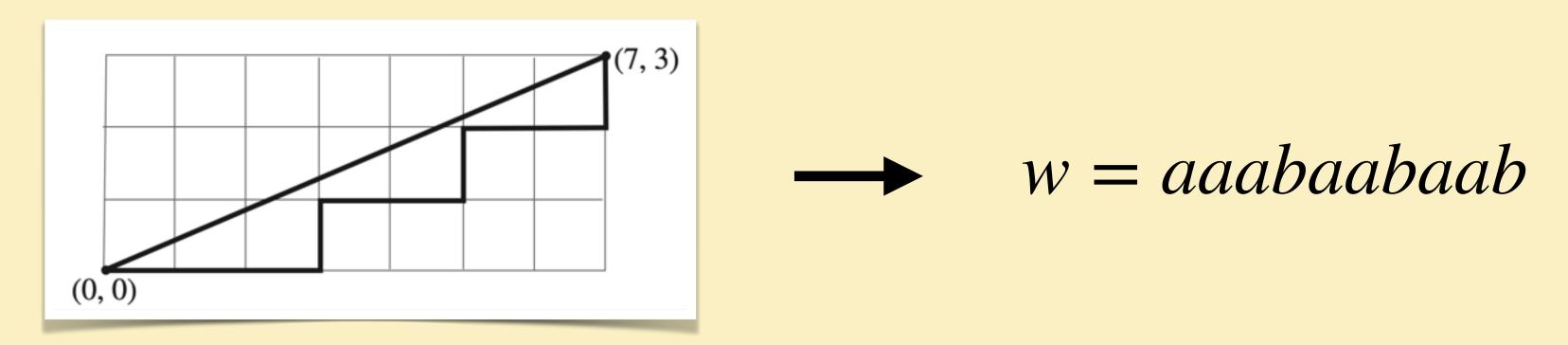
Consequently, such a word cannot be Lyndon.

Corollary. All $\mathscr C$ -long Lyndon words are k-perfect.

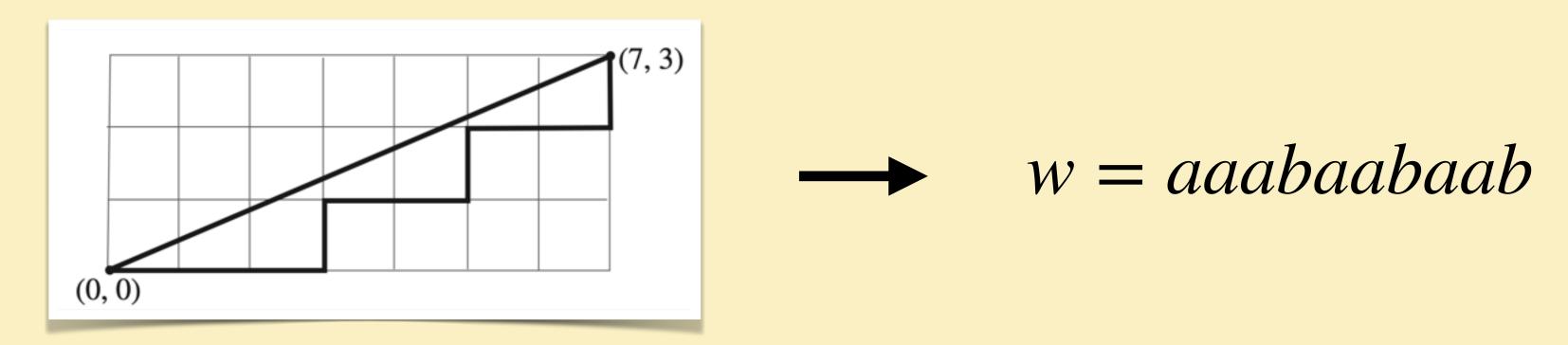
$$\pi_{\sigma}(k, \ell = k+1) = \lambda_{\sigma}(\ell)$$

Explicit formula exists

Definition (Christoffel word). A word on a two-letter alphabet is a (lower) Christoffel word if it is obtained by discretizing a segment in the plane.

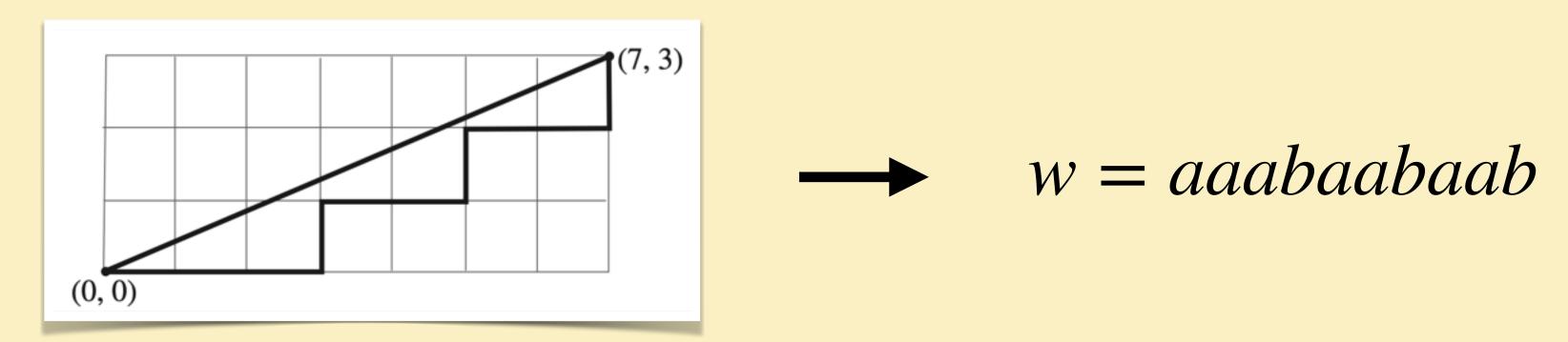


Definition (Christoffel word). A word on a two-letter alphabet is a (lower) Christoffel word if it is obtained by discretizing a segment in the plane.



It is known that a word of length ℓ that has $(\ell-1)$ distincts factors of length $(\ell-2)$ is a rotation of a Christoffel word.

Definition (Christoffel word). A word on a two-letter alphabet is a (lower) Christoffel word if it is obtained by discretizing a segment in the plane.



It is known that a word of length ℓ that has $(\ell-1)$ distincts factors of length $(\ell-2)$ is a rotation of a Christoffel word.

$$\pi_{\sigma}(k, \ell = k + 2) = \lambda_{\sigma}(\ell) - \varphi(\ell) \cdot \begin{pmatrix} \sigma \\ 2 \end{pmatrix}$$

(number of Christoffel words) x (choice of two letters)

Conjecturing $\pi_{\sigma}(k,\ell=k+3)$

Let
$$\psi$$
 be the function: $n \mapsto 3/2 \cdot \varphi(n)$ if $n \equiv 0 \mod 4$, $n \mapsto 2 \cdot \varphi(n)$ if $n \equiv 2 \mod 4$, $n \mapsto \varphi(n)$ otherwise.

$$\pi_{\sigma}(k,\ell=k+3) = \lambda_{\sigma}(k,\ell) - \binom{\sigma}{2} (\sigma \cdot \psi(\ell) - 2) \tag{?}$$

Conjecturing $\pi_{\sigma}(k,\ell=k+3)$

Let
$$\psi$$
 be the function: $n \mapsto 3/2 \cdot \varphi(n)$ if $n \equiv 0 \mod 4$, $n \mapsto 2 \cdot \varphi(n)$ if $n \equiv 2 \mod 4$, $n \mapsto \varphi(n)$ otherwise.

$$\pi_{\sigma}(k, \ell = k+3) = \lambda_{\sigma}(k, \ell) - \binom{\sigma}{2} (\sigma \cdot \psi(\ell) - 2) \tag{?}$$

Process.

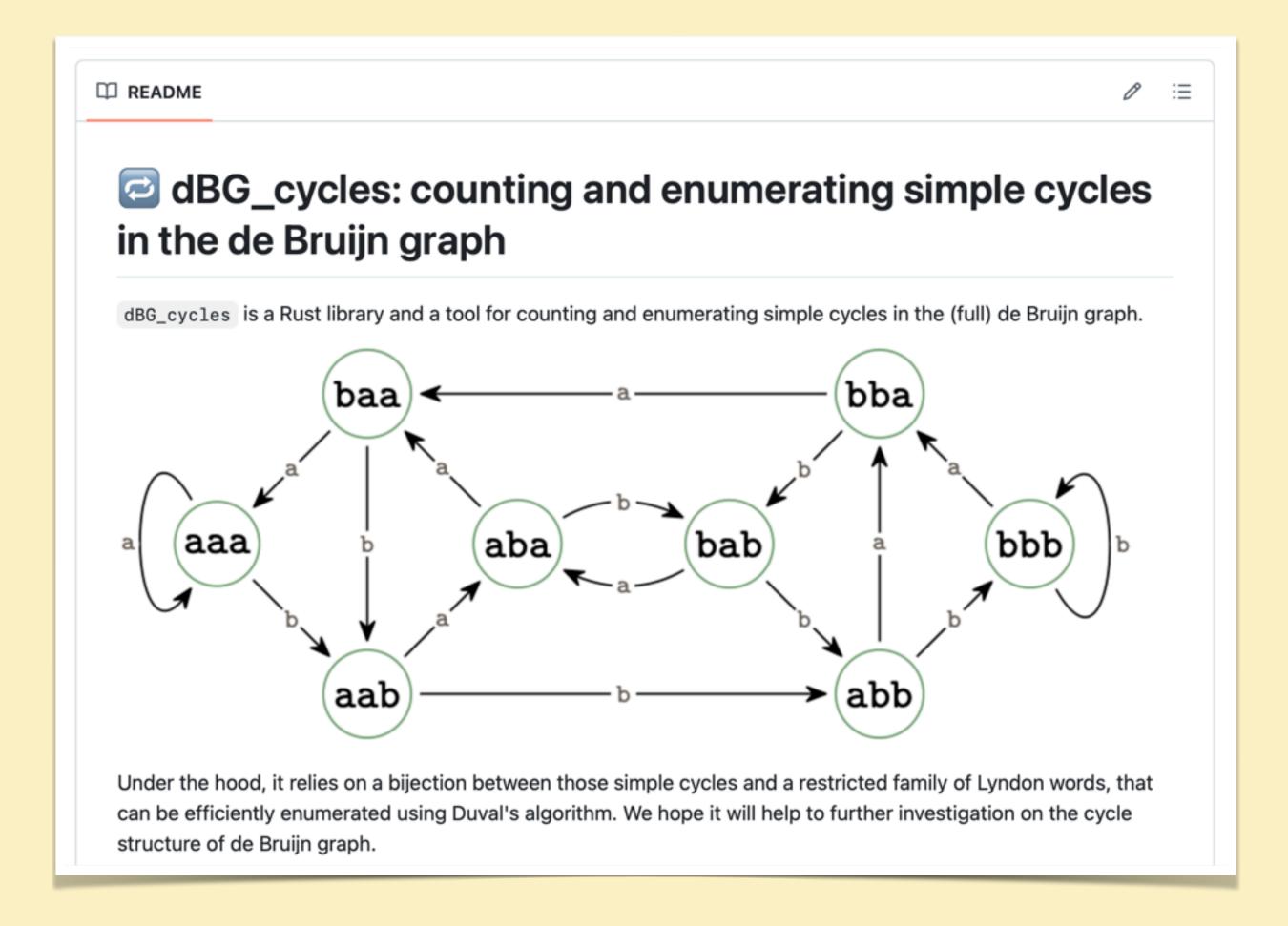
- 1. We found out that $(\bar{\pi}_2(k, l = k + 3))_k = 2 \cdot \text{OEIS}_{A126246} 2$ using OEIS,
- 2. Guessing a small impact of σ on the formula, we found that $\bar{\pi}_3 = 6 \cdot \text{OEIS}_{\text{A126246}} 9$, $\bar{\pi}_4 = 24 \cdot \text{OEIS}_{\text{A126246}} 12$, $\bar{\pi}_5 = 50 \cdot \text{OEIS}_{\text{A126246}} 20$ and $\bar{\pi}_6 = 90 \cdot \text{OEIS}_{\text{A126246}} 30$,

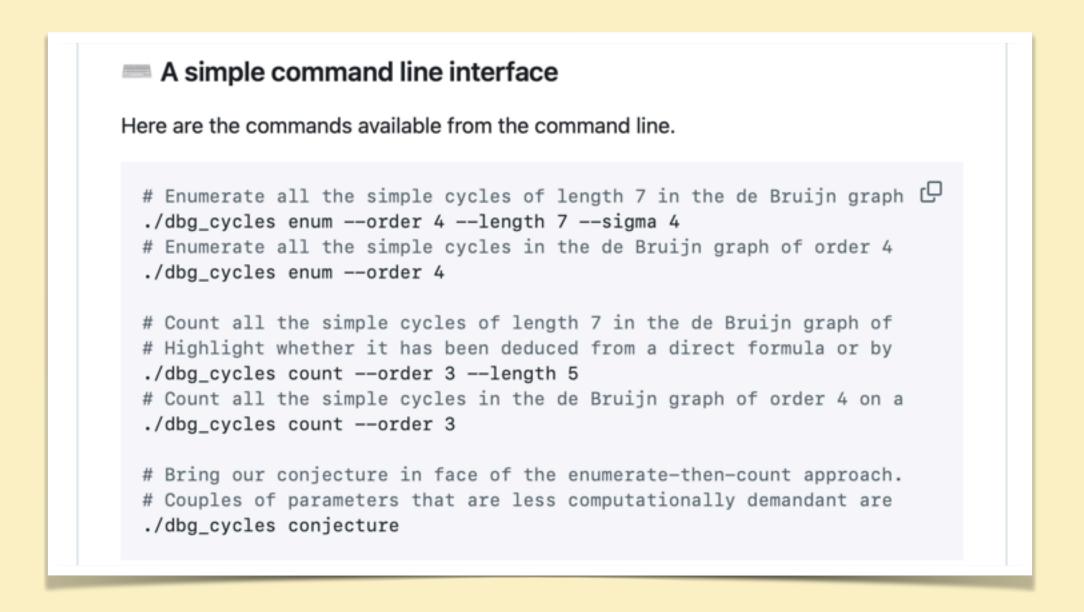
3. We derive the formula, and verify it on tractable instances.

A small Rust library

The dbg cycles library

Objective. Make this counting as accessible as possible.

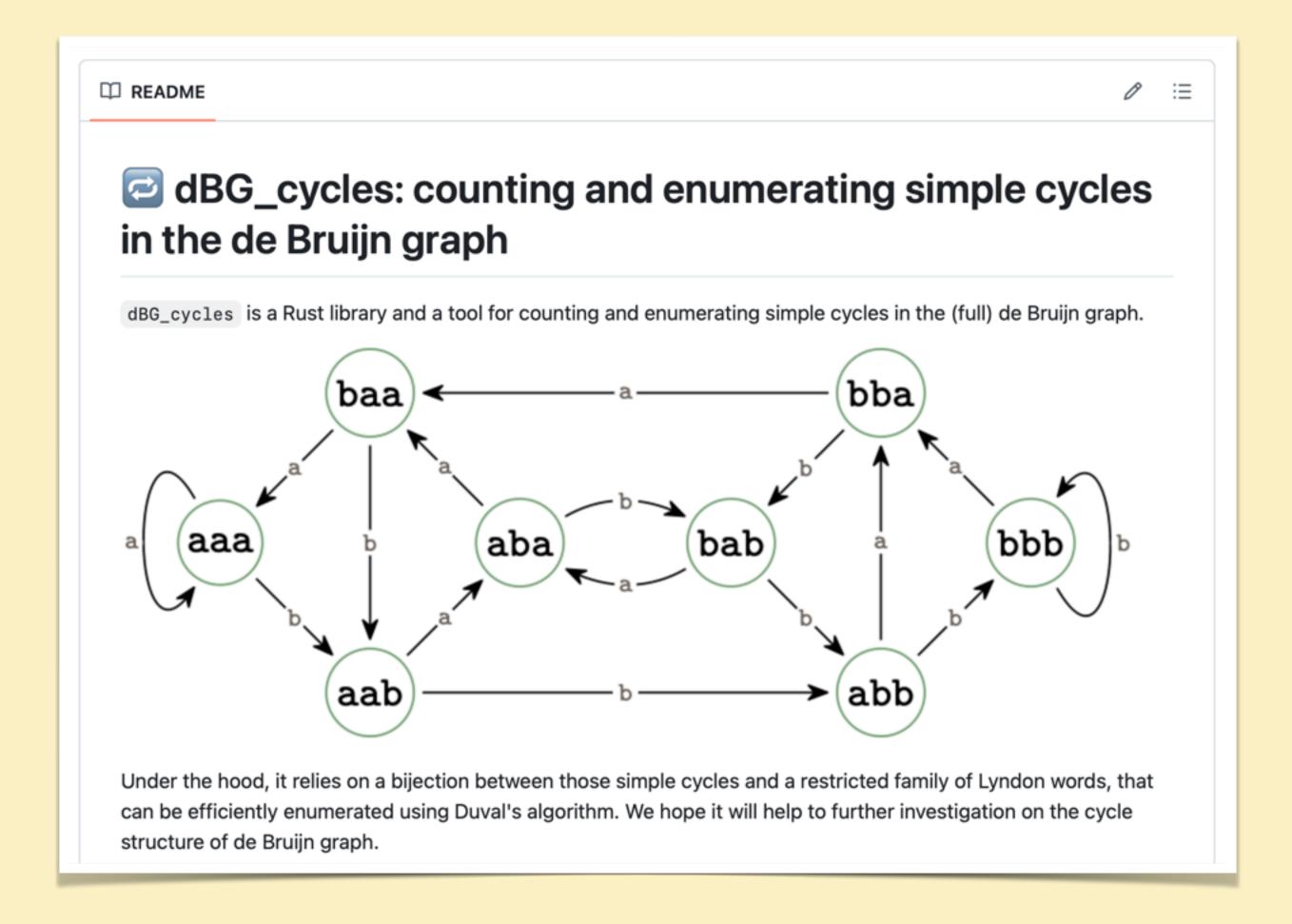


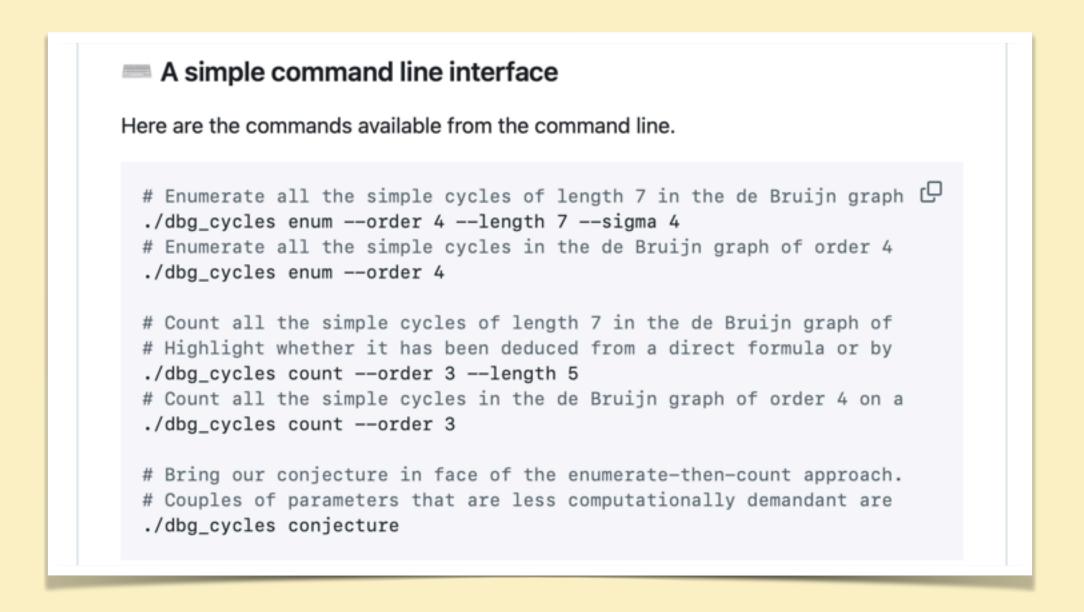


It implements the Duval's algorithms (Lyndon word enumeration) and the bijection between Lyndon words and dBG simple cycles.

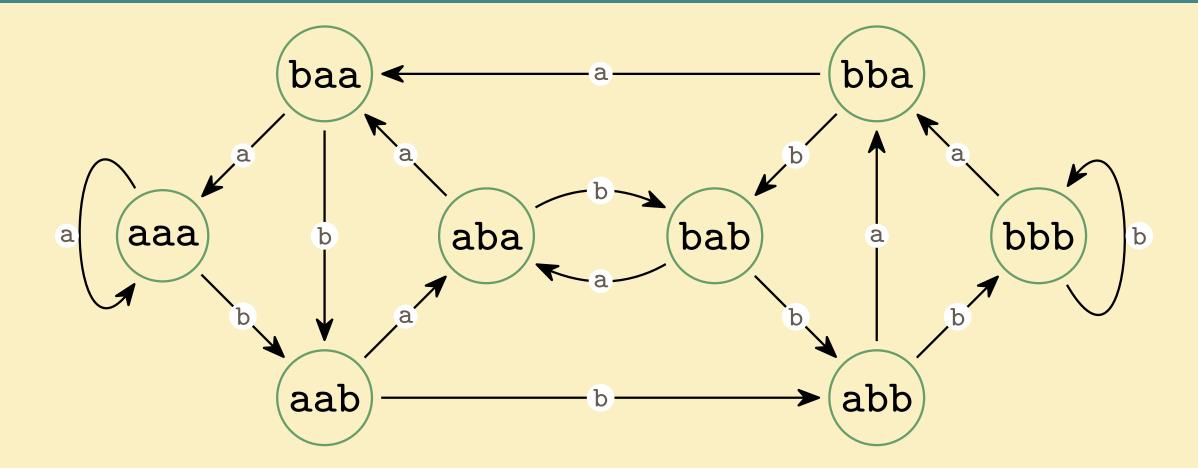
The dbg cycles library

Objective. Make this counting as accessible as possible.





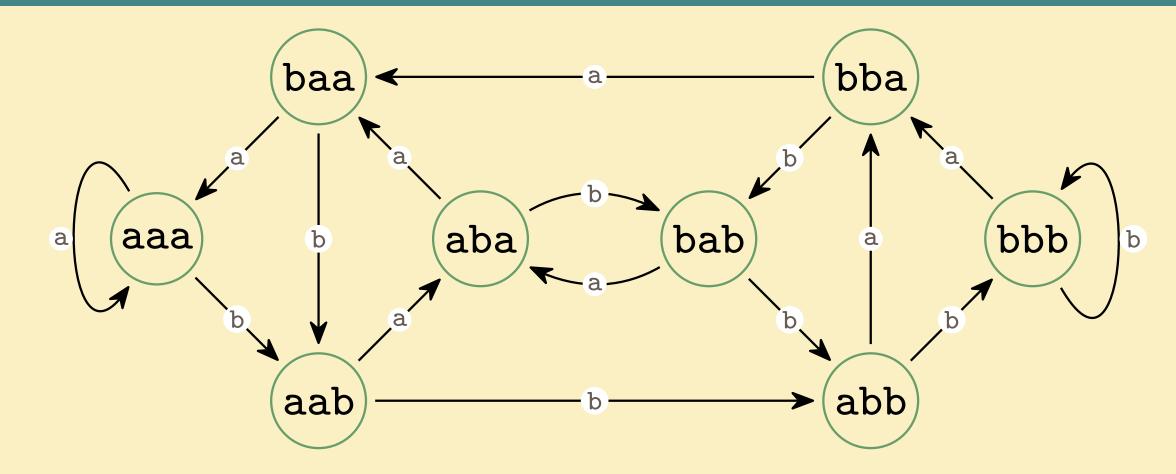
It implements the Duval's algorithms (Lyndon word enumeration) and the bijection between Lyndon words and dBG simple cycles. **Limitation.** very slow for enumerating cycles larger than the order of the graph (filtering step, that discards most candidates).



Interest for bioinformatics. Sliding windows (minimizers, compression schemes)

Their topology is still not fully understood. Eg. their simple cycles

- we collected the state of knowledge
- we slightly extended it
- we made these results as accessible as possible

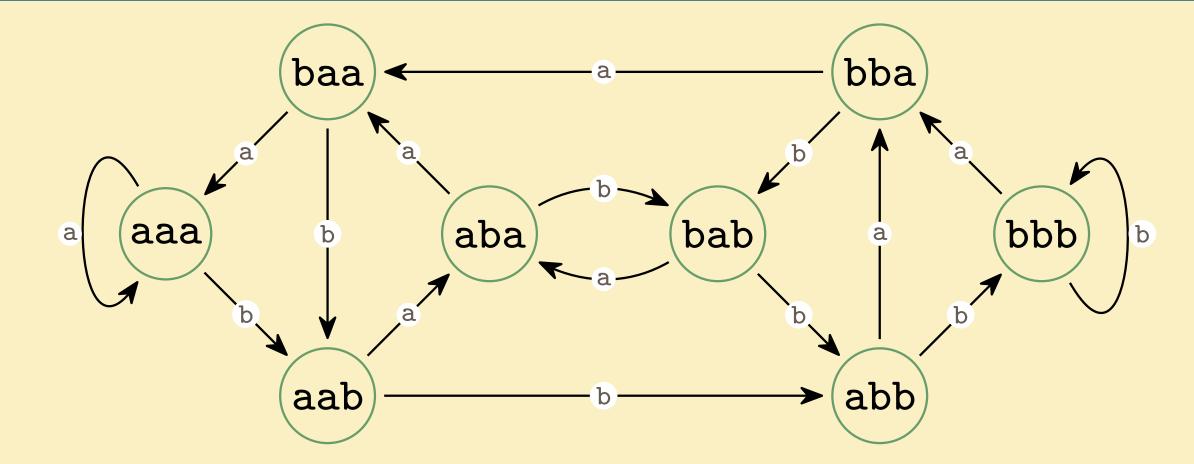


Interest for bioinformatics. Sliding windows (minimizers, compression schemes)

Their topology is still not fully understood. Eg. their simple cycles

- we collected the state of knowledge
- we slightly extended it
- we made these results as accessible as possible

Open question 1. Can we find (and prove) explicit formulas for the number of simple cycles for the regimes $\ell \ge k+3$?



Interest for bioinformatics. Sliding windows (minimizers, compression schemes)

Their topology is still not fully understood. Eg. their simple cycles

- we collected the state of knowledge
- we slightly extended it
- we made these results as accessible as possible

Open question 1. Can we find (and prove) explicit formulas for the number of simple cycles for the regimes $\ell \ge k + 3$?

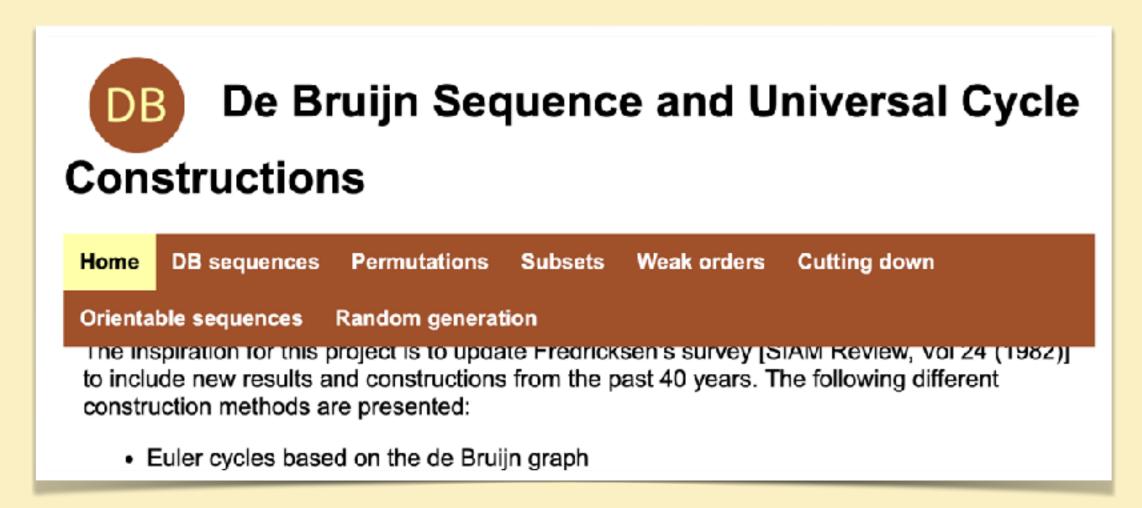
Open question 2. Can we adapt Duval's algorithms to enumerate k-perfect Lyndon words directly? in optimal space and time?

Thoughts

- History of dBG and Lyndon words is long ⇒ hard to navigate between equivalent/close notions and "folklore" results when non-specialist.
- Some papers are in French (eg. Duval's Algorithms): how accessible are these results for the international community?

Thoughts

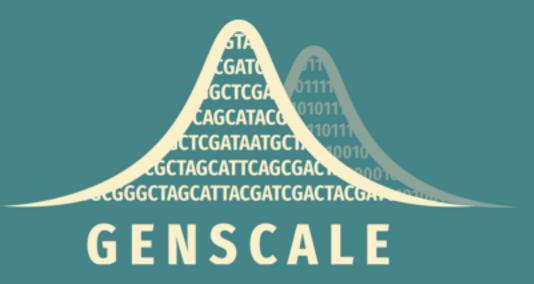
- History of dBG and Lyndon words is long ⇒ hard to navigate between equivalent/close notions and "folklore" results when non-specialist.
- Some papers are in French (eg. Duval's Algorithms): how accessible are these results for the international community?
- Some nice encounters:





Open question 3. Should we start a collective effort and launch debruijngraph.org?





Counting simple cycles

in the de Bruijn graph

Léo Ackermann, Pierre Peterlongo

